# Proxy-bridged Image Reconstruction Network for Anomaly Detection in Medical Images

Kang Zhou, Jing Li, Weixin Luo, Zhengxin Li, Jianlong Yang,
Huazhu Fu, Jun Cheng, Jiang Liu and Shenghua Gao

*Abstract*—Anomaly detection in medical images refers to the identification of abnormal images with only normal images in the training set. Most existing methods solve this problem with a self-reconstruction framework, which tends to learn an identity mapping and reduces the sensitivity to anomalies. To mitigate this problem, in this paper, we propose a novel Proxy-bridged Image Reconstruction Network (ProxyAno) for anomaly detection in medical images. Specifically, we use an intermediate proxy to bridge the input image and the reconstructed image. We study different proxy types, and we find that the superpixel-image (SI) is the best one. We set all pixels' intensities within each superpixel as their average intensity, and denote this image as SI. The proposed ProxyAno consists of two modules, a Proxy Extraction Module and an Image Reconstruction Module. In the Proxy Extraction Module, a memory is introduced to memorize the feature correspondence for normal image to its corresponding SI, while the memorized correspondence does not apply to the abnormal images, which leads to the information loss for abnormal image and facilitates the anomaly detection. In the Image Reconstruction Module, we map an SI to its reconstructed image. Further, we crop a patch from the image and paste it on the normal SI to mimic the anomalies, and enforce the network to reconstruct the normal image even with the pseudo abnormal SI. In this way, our network enlarges the reconstruction error for anomalies. Extensive experiments on brain MR images, retinal OCT images and retinal fundus images verify the effectiveness of our method for both image-level and pixel-level anomaly detection.

*Index Terms*—Anomaly Detection, Proxy, Superpixel-Image, Memory, Pseudo Anomalies

## I. INTRODUCTION

Anomaly detection aims to identify abnormalities with only normal data in the training set [1], [2], [3]. Recently it has drawn much attention in the community of medical image analysis [4], [5], [6], [7]. There are three reasons for this: firstly, it is not easy to acquire medical data, especially for some rare diseases; secondly, it is expensive to annotate the lesion; thirdly, it is comparatively easier to collect the data from healthy subjects. Because of the data constraint, anomaly detection in medical images is more challenging than disease classification that has been well tackled with deep convolutional neural networks (CNNs) [8], [9], [10], [11].

To tackle the anomaly detection with only normal data, many methods have been proposed. These methods can be roughly categorized into two groups. The first group is reconstruction-based method, which train a model to reconstruct the normal image, and in the test phase, the reconstruction error of abnormal images is larger than the normal ones. The second group is non-reconstruction-based methods. For example, Ouardini *et al.* [12] proposed an efficient and effective transfer-learning based approach for anomaly detection on retinal fundus images. Golan *et al.* [13] proposed to learn a meaningful representation of the learned training data in a fully discriminative fashion using the self-labeled dataset.

In this work, we focus on the study of reconstruction-based anomaly detection. Most previous reconstruction-based methods follow a self-reconstruction paradigm. Specifically, an Auto-Encoder (AE) is commonly used for anomaly detection [14], [15]. In the training phase, the AE is learnt with only normal data, and it is expected that the reconstructed image from the decoder will be close to the input for the normal image, and the output is different from the input for the abnormal images in the test phase. In addition, generative adversarial network (GAN) [16] based approaches have also be introduced to guarantee the fidelity of the reconstruction for normal data [17], [4]. Further, other than measuring the anomaly with reconstruction error in the image space, Akcay *et al.* [18] proposed to append another encoder to extract latent features corresponding to the reconstructed image, and take the difference of latent features corresponding to the input and reconstructed image as the anomaly measurement. However, it is worth noting that self-reconstruction is essentially to learn an identity mapping function due to the information equivalence, *i.e.*, the equivalence between the input and output of the model during the training [19]. Specifically, self-reconstruction aims to reproduce the output to approximate the input. They may hence tend to overfit to learn an identity mapping between the input and output [20]. Therefore, self-reconstruction-based anomaly detection methods cannot guarantee the large reconstruction error for abnormal images, such that it may not be

Corresponding Author: Shenghua Gao, gaoshh@shanghaitech.edu.cn.

Kang Zhou, Jing Li, Weixin Luo, Zhengxin Li and Shenghua Gao are with School of Information Science and Technology, ShanghaiTech University, Shanghai 201210, China; Kang Zhou and Jing Li are also with Shanghai Institute of Microsystem and Information Technology, Chinese Academy of Sciences, Shanghai 200050, China, and also University of Chinese Academy of Sciences, China. Kang Zhou and Jing Li contributed equally to this work. Email: {zhoukang, lijing1}@shanghaitech.edu.cn. Shenghua Gao is also with Shanghai Engineering Research Center of Intelligent Vision and Imaging, Shanghai 201210, China, and also with Shanghai Engineering Research Center of Energy Efficient and Custom AI IC, Shanghai 201210, China.

Jianlong Yang is with School of Biomedical Engineering, Shanghai Jiao Tong University, Shanghai 200030, China.

Huazhu Fu is with the Institute of High Performance Computing (IHPC), Agency for Science, Technology and Research (A*STAR), Singapore 138632.

Jun Cheng is with the Institute for Infocomm Research, Agency for Science, Technology and Research (A*STAR), Singapore 138632.

Jiang Liu is with the Department of Computer Science and Engineering, Southern University of Science and Technology, Guangdong 518055, China, and also with Cixi Institute of Biomedical Engineering, Chinese Academy of Sciences, Zhejiang 315201, China.

sensitive to anomalies.

To mitigate the identity mapping problem in self-reconstruction-based anomaly detection, we use an intermediate proxy to bridge the input image and the reconstructed image. We study different proxy types (*e.g.*, edge, smooth image), and we find that the superpixel-image (SI) is the best one. As shown in Fig. 1 (A), we set all pixels' intensities within each superpixel to their average intensity for each superpixel and denote this image as SI. In our proposed anomaly detection method, we map the input image to a proxy with the Proxy Extraction Module, and map the proxy to its reconstructed image with the Image Reconstruction Module. We name our method as Proxy-bridged Image Reconstruction Network (ProxyAno).

**The motivations behind using SI as the proxy are the following aspects. 1) From the view of biomedical analogy.** First of all, we would like to argue that the SI is an analogy to the tissue. Concretely, the scale of tissue is between that of the cell and the organ, and the tissue is a group of many similar cells and carries out a specific function. Similarly, the superpixel is defined as a group of pixels with similar characteristics (*e.g.*, intensity) and the scale of superpixel is between that of the pixel and the image [21]. As shown in Fig. 1 (B), we take the eyeball as the specific example to illustrate the analogy. When a disease occurs, the cells and tissues in the lesion region undergo pathological changes [22], [23], [24]. In other words, the tissue-level pathological change represents the disease occurring. Analogously, the SI-level change partly represents the anomaly that occurs in the given image [25]. This analogy inspires us to use the SI as an intermediate proxy for image reconstruction. **2) From the view of frequency domain.** In medical images, the average gradient (*i.e.*, frequency) of the abnormal region is usually larger than that of the normal region [26], [27]. By transforming the original image to its SI, the abnormal region loses more information than the normal region. Thus, the variation degree of the abnormal region is more significant than that of the normal region. Therefore, by leveraging SI as the intermediate proxy, the reconstruction from SI of the abnormal samples is more difficult than that of the normal samples, leading to favoring detecting anomalies. We take a specific example in Fig. 1 (C). The images in the brain MRI dataset [28] are transferred with 2D Fourier Transform [29] and the log-spectral amplitude is estimated. It can be observed that the major distinction between normal images and abnormal images is the high-frequency information, and the decrease degree of the high-frequency's amplitude of the abnormal samples is more significant than that of the normal samples. **3) From the definition of superpixel.** The definition of superpixel determines that the SI can preserve both the local structure and texture information [30], [31]. In other words, SI contains mid-level semantic information [32].

By using the SI as the intermediate proxy to bridge the input image and the reconstructed image, it is expected that the network will lead to a large prediction error in SI for an abnormal image, and consequently will lead to a large reconstruction error for an abnormal image. The proposed ProxyAno consists of two modules, *i.e.*, a Proxy Extraction Module to map the input image to an SI, and an Image Reconstruction Module to map the SI to the reconstructed image. Specifically, in the Proxy Extraction Module, we introduce a memory to memorize the mapping pattern from the input image to its corresponding SI for normal training data. Given an input, we use its latent feature extracted by the encoder in the Proxy Extraction Module to retrieve the most similar item in the memory, and feed the retrieved feature into the decoder in Proxy Extraction Module to predict the SI. Then we feed the predicted SI into the Image Reconstruction Module to reconstruct the input image. To further enlarge the reconstruction error for abnormal images, we create pseudo abnormal SI by cutting a patch from a randomly selected normal image and pasting the patch on the normal proxy. Other than the reconstruction from the normal SI, we also enforce the network well reconstruct the input even with the pseudo abnormal SI. In this way, the Image Reconstruction Module would have the ability to repair the anomaly, which would result in a large reconstruction error for abnormal images. Experiments on the images of different modalities validate the effectiveness of our approach for both image-level and pixel-level anomaly detection.

The main contributions are summarized as follows:

1) To mitigate the identity mapping problem in self-reconstruction-based anomaly detection, under the assumption that the mapping from an image and its SI can be predictable for the normal images while the mapping for abnormal images cannot be predicted, we propose the ProxyAno to connect the input image and the reconstructed image with an SI.

2) In the Proxy Extraction Module, we introduce a memory to memorize the correspondence between the input image and its SI for normal images, while the memorized correspondence does not apply to the abnormal images.

3) In the Image Reconstruction Module, a pseudo abnormal SI is built to enforce the network to reconstruct the normal input even with the abnormal SI. This would further enlarge the reconstruction error for abnormal image.

4) Extensive experiments on three modalities (*i.e.*, brain MRI, retinal OCT and retinal fundus) demonstrate the effectiveness of our approach for anomaly detection. Moreover, our work verifies the effectiveness of anomaly detection for abnormal region segmentation task in retinal fundus images.

The rest of this paper is organized as follows: In Section II, we introduce the work related to anomaly detection in medical images. In Section III, we describe our proposed ProxyAno for anomaly detection in detail. In Section IV, extensive experiments on three modalities are conducted to validate the effectiveness of our method. We conclude our work in Section V.

## II. RELATED WORK

Recently, anomaly detection has drawn much attention in the medical image domain, including different imaging modalities, such as brain MRI [5], [6], [7], [14], [33], [34],
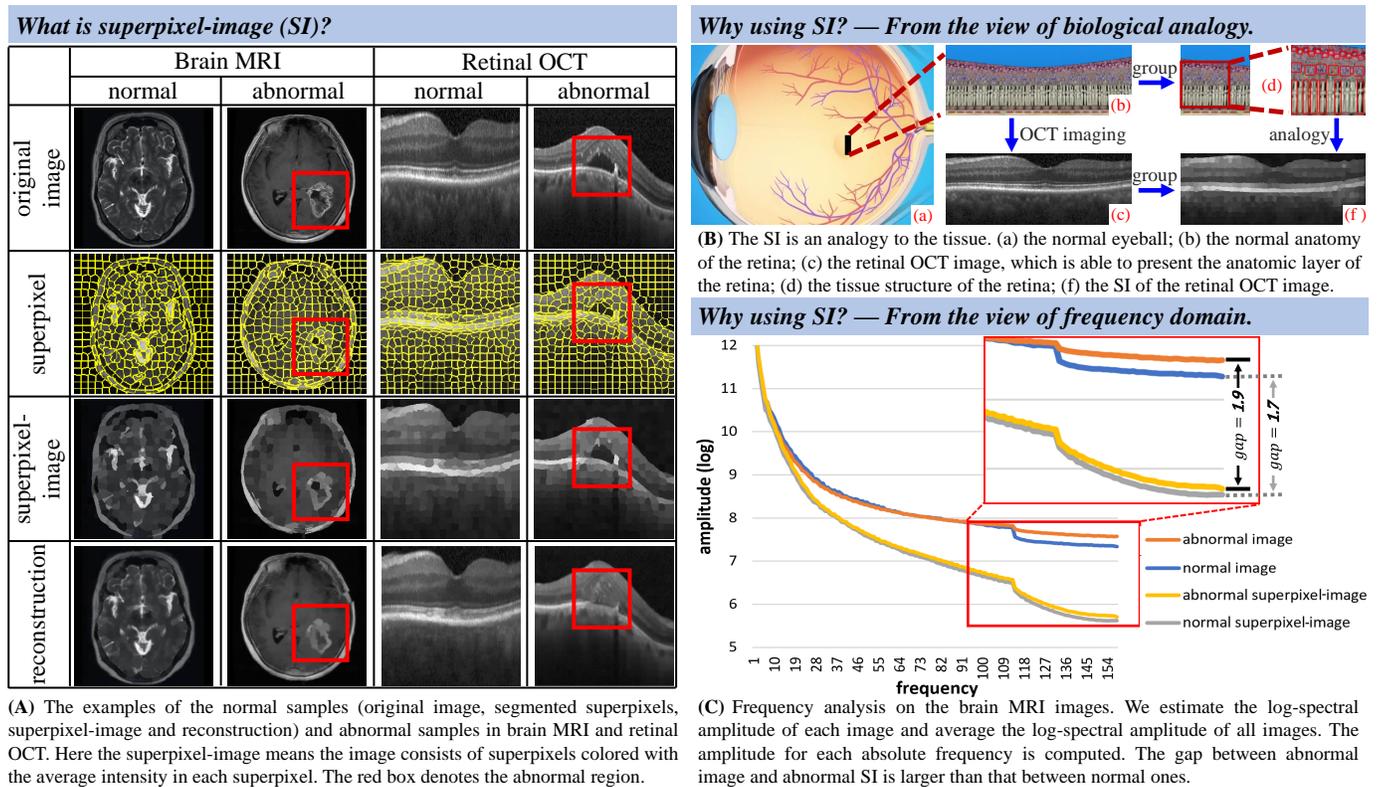
**What is superpixel-image (SI)?**

**(A)** The examples of the normal samples (original image, segmented superpixels, superpixel-image and reconstruction) and abnormal samples in brain MRI and retinal OCT. Here the superpixel-image means the image consists of superpixels colored with the average intensity in each superpixel. The red box denotes the abnormal region.

**Why using SI? — From the view of biological analogy.**

**(B)** The SI is an analogy to the tissue. (a) the normal eyeball; (b) the normal anatomy of the retina; (c) the retinal OCT image, which is able to present the anatomic layer of the retina; (d) the tissue structure of the retina; (f) the SI of the retinal OCT image.

**Why using SI? — From the view of frequency domain.**

**(C)** Frequency analysis on the brain MRI images. We estimate the log-spectral amplitude of each image and average the log-spectral amplitude of all images. The amplitude for each absolute frequency is computed. The gap between abnormal image and abnormal SI is larger than that between normal ones.

Fig. 1. The definition of SI and the motivations of using SI as the intermediate proxy to bridge the input image and the reconstructed image.

retinal OCT [4], [35], [36], [37], chest X-Ray [38], [39], [40], [41] and brain CT [42]. In this paper, we focus on the reconstruction-based anomaly detection in medical images. These works can be roughly categorized into two categories: image-level anomaly detection and pixel-level anomaly detection. The pixel-level anomaly detection also called as abnormal region segmentation and anomaly segmentation, which classifies the pixels in the given image into normal or abnormal.

**Image-Level Anomaly Detection.** Most previous image-level anomaly detection methods are based on a self-reconstruction framework that assumes the abnormal images cannot be well reconstructed by a model learned merely with the normal images. Specifically, Zimmerer *et al.* proposed to use a variational auto-encoder for the brain MRI image reconstruction [5], and they also proposed to combine a context-encoder and a variational auto-encoder for the brain MRI image reconstruction [33]. Schlegl *et al.* [4] proposed to utilized a GAN [16] for image reconstruction. They first trained a generator mapping the latent vector to the reconstructed image. Then they fixed the generator and the discriminator to train a mapping from the input image to a latent vector with the back-propagation algorithm. Based on [4], Han *et al.* [34] proposed to take multiple adjacent slices of 3D MRI data as the input for the reconstruction. Further, to address the slow mapping issue from an image to a latent vector in [4], Schlegl *et al.* [35] proposed to use an encoder to map the input image to the latent vector. Inspired by the success of pix2pix [43], which trains an encoder and a generator with an end-to-end manner, Zhou *et al.* [36] proposed to use the

pix2pix and regularize the sparsity of the latent space (termed Sparse-GAN) to guarantee the fidelity of the reconstruction for normal data. Ouardini *et al.* [12] proposed to employ the Inception-ResNet-v2 network [44] trained on ImageNet dataset [45] as a feature extractor, and subsequently feeds the extracted feature into Isolation Forests [46], which is a robust and efficient anomaly detection method. They applied their method to image-level anomaly detection on retinal fundus image and conducted experiments on two datasets, which contain Retinopathy of Prematurity (ROP) and Diabetic Retinopathy (DR), respectively. The model [12] can classify the given image into normal or abnormal (ROP and DR). Very recently, to leverage the structure information (*e.g.*, the retinal blood vessels in fundus image, and the retinal layers in OCT image) for anomaly detection, Zhou *et al.* [47] proposed a P-Net to reconstruct the image, which leverages the relation between structure and texture for image-level anomaly detection. Zhou *et al.* [47] apply image-level anomaly detection to multiple retinal disease detection. Since image-level anomaly detection is a two-classes classification problem, they regard the images with DR, glaucoma, age-related macular degeneration and pathological myopia as the abnormal images. Then the P-Net [47] classifies the given image into normal or abnormal.

**Pixel-Level Anomaly Detection.** Besides used for image-level anomaly detection, reconstruction based methods are also commonly used for pixel-level anomaly detection. Most existing methods train a reconstruction model with the normal data, and in the test phase, the residual map (a.k.a anomaly map, lesion map) is obtained by subtract reconstructed image with the input image. Chen *et al.* [17] initially proposed to use
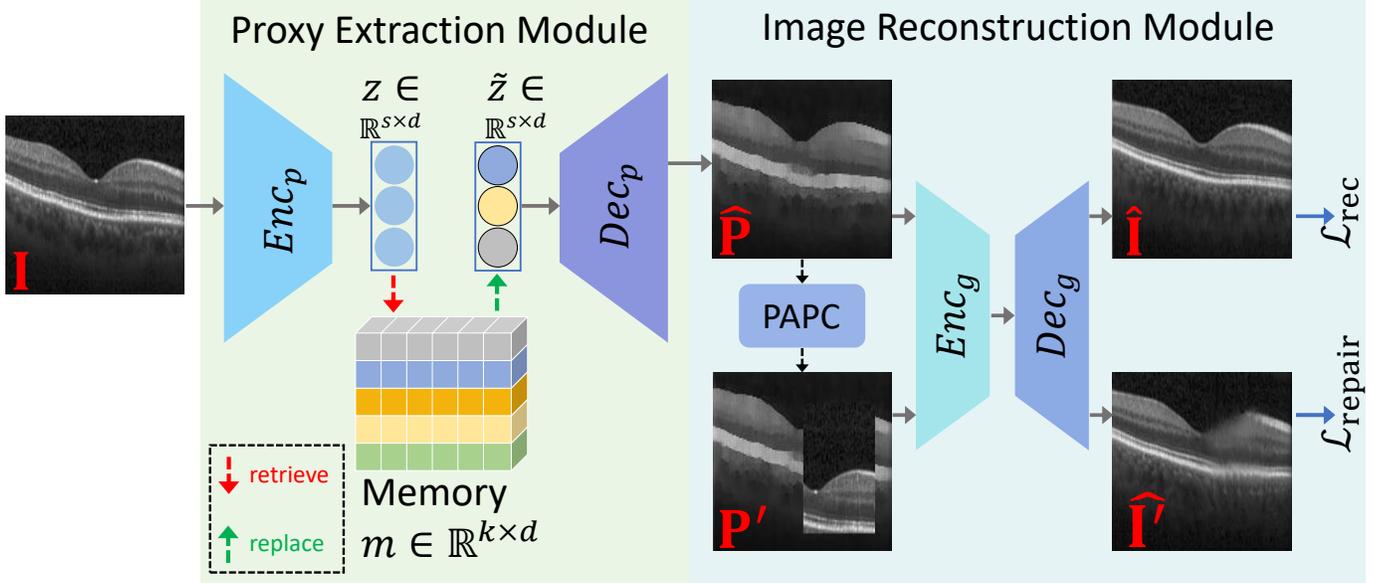
Fig. 2. The overview of the proposed network, which consists of two modules, a Proxy Extraction Module and an Image Reconstruction Module. PAPC is the short for the Pseudo Abnormal Proxy Constructor, and PAPC only works in the training phase. $\hat{\mathbf{P}}$ and $\mathbf{P}'$ denote the predicted proxy and pseudo abnormal proxy, respectively.

adversarial auto-encoders for anomaly segmentation on brain MRI images. Almost the same time, Baur *et al.* [14] proposed to use a deep auto-encoder that combines spatial AEs and GANs for anomaly segmentation on brain MRI. In addition, the Bayesian deep learning also been introduced for the pixel-level anomaly detection. Pawlowski *et al.* [42] proposed to use the Bayesian auto-encoder to model the normal data distribution and apply the algorithm for anomaly detection on brain CT. Seebock *et al.* [37] introduced a Bayesian U-Net that exploits the segmentation model of normal anatomy and its epistemic uncertainty for anomaly segmentation on the test images. Very recently, Chen *et al.* [7] proposed a probabilistic model that uses a network-based prior as the normative distribution and detects lesions using Maximum-A-Posteriori estimation. Baur *et al.* [6] further proposed a dropout-U-Net based auto-encoder, which introduces the dropout connections between the encoder and decoder, to model the uncertainty and enable high fidelity reconstructions on brain MRI images.

**Summary.** In a summary, almost all previous works adopted the self-reconstruction paradigm for anomaly detection in medical images. As aforementioned, the self-reconstruction model is essentially to learn an identity mapping function, which cannot guarantee a large reconstruction error for abnormal images. To mitigate the identity mapping problem, we propose to introduce an SI as the proxy to bridge the input image and the reconstructed image. Moreover, different with [47], [12], which conducted image-level anomaly detection (*i.e.*, classifying an image into normal or abnormal) on the retinal fundus image, this work verifies the effectiveness of pixel-level anomaly detection on retinal fundus images. To the best of our knowledge, we are the first to conduct anomaly segmentation in the retinal fundus modality.

## III. METHOD

### A. Overview

In the reconstruction based anomaly detection framework, it is expected that the network can well reconstruct the normal input image, and the reconstructed image of the abnormal ones would be with a large reconstruction error. In this way, the abnormal images can be distinguished from the normal images. Towards this end, in this paper, we introduce a Proxy-bridged Image Reconstruction Network (ProxyAno) for anomaly detection in medical images. The proposed ProxyAno consists of two modules, a Proxy Extraction Module $\mathcal{F}_p(\cdot)$ and an Image Reconstruction Module $\mathcal{F}_g(\cdot)$. Specifically, the proposed ProxyAno can be formulated as:

$$\begin{aligned} \hat{\mathbf{P}} &= \mathcal{F}_p(\mathbf{I}) \\ \hat{\mathbf{I}} &= \mathcal{F}_g(\hat{\mathbf{P}}) \end{aligned} \tag{1}$$

where $\hat{\mathbf{P}}$ denotes the predicted SI, $\mathbf{I}$ and $\hat{\mathbf{I}}$ denote original image and reconstructed image, respectively. We use $\boldsymbol{Enc_p}/\boldsymbol{Dec_p}$ to denote the encoder/decoder in the Proxy Extraction Module, and use $\boldsymbol{Enc_g}/\boldsymbol{Dec_g}$ to denote the encoder/decoder in the Image Reconstruction Module, respectively.

Fig. 2 illustrates the overview of our ProxyAno. To enlarge the reconstruction error for the abnormal image while maintaining a small reconstruction error for the normal images, we propose to: 1) memorize the mapping pattern between the normal input image and its corresponding SI with a memory in the Proxy Extraction Module. Once this module is well trained, the weights of this module are fixed when training the Image Reconstruction Module; 2) reconstruct the image from its SI, meanwhile repairing the abnormal SI in the Image Reconstruction Module. To achieve this, we first create a pseudo abnormal proxy $\mathbf{P}'$ and feed it into the Image Reconstruction Module with $\hat{\mathbf{I}}' = \mathcal{F}_g(\mathbf{P}')$. Then we propose a repairing loss to enforce the similarity between the

reconstructed image $\hat{\mathbf{I}}'$ of the pseudo abnormal proxy and the normal image $\mathbf{I}$. Benefiting from these two improvements, in the test phase, both $\hat{\mathbf{P}}$ and $\hat{\mathbf{I}}$ tend to be close to the normal ones, leading to a small reconstruction error for the normal samples and a large reconstruction error on the abnormal ones. Then the reconstruction error is used as a measurement to detect the anomalies.

### B. Superpixel-image-bridged Image Reconstruction Network

In our work, we use the SLIC [48] algorithm to obtain the superpixels. The SLIC algorithm groups meaningful pixels into a superpixel by composing spatially adjacent pixels. Based on the superpixels, we can get an SI that consists of the superpixel colored with the average intensity for pixels within each superpixel. Then we use the SI extracted with the SLIC algorithm as supervision to train a Proxy Extraction Module. It is worth noting that directly applying the SLIC algorithm on the abnormal image, the superpixel structure of lesion can also be extracted. Thus, the abnormal SI contains some lesion information; consequently, the lesion in the abnormal images may also be reconstructed. To increase the fidelity of the reconstruction of the normal images and increase the reconstruction error for the abnormal ones, on the one hand, we introduce a memory to map the abnormal image to a normal SI, on the other hand, we also propose to let the algorithm automatically repairing an SI of an abnormal image, and map it to a normal input image.

### C. Proxy Extraction Module

In the Proxy Extraction Module, we propose to augment the network with a memory to explicitly memorize the correspondence between the normal input and its SI. The motivations behind using a memory to augment the network are three aspects: 1) to recognize disease in medical images, the clinicians need to memorize the characteristics of normal samples, and the superpixel structure is an important characteristic in the image; 2) the memory is a type of location-independent long-term attention, which augments the ability of convolutional structure; 3) the correspondence between the normal input and its SI is very regular. Using the correspondence patterns is probably enough to generate all normal SI's from its input. Therefore, we introduce a memory that can memorize the normal patterns. The feature is first extracted with a encoder. Then, instead directly feed the feature into the decoder, we use the feature as a query to retrieve the most relevant item in the memory. The features feed into the decoder are obtained from a selected memory item of the normal data.

Since this memory is learnt based on the correspondence of normal images, it can be expected that this would cause information loss for the abnormal ones when generating its SI. Consequently, the image reconstruction form this SI would be with a large reconstruction error, which is a desirable property for anomaly detection. To achieve this goal, we introduce a memory $m \in \mathbb{R}^{k \times d}$ to memorize the latent features in the Proxy Extraction Module, where $k$ is the size (the number of items) of the memory , and $d$ is the dimensionality of the latent vector $m_j \in \mathbb{R}^d, j \in \{1, 2, \cdots, k\}$. As shown in Fig.

2, the framework takes an image $\mathbf{I}$ as an input, and then $\mathbf{I}$ is passed through the encoder $\boldsymbol{Enc}_p$ to extract a latent feature $z \in \mathbb{R}^{h \times w \times d}$, where $h$ and $w$ are the height and the width of the feature map. In this module, $z$ is used to search for the nearest neighbor item in the memory $m$, and we denote the nearest item of $z$ as $\tilde{z} \in \mathbb{R}^{s \times d}$, here $\tilde{z}$ indicates the nearest item in the memory of a feature.

Specifically, as illustrated in Fig. 3, we first flatten the latent feature as $z \in \mathbb{R}^{s \times d}$, where $s = h \times w$. Then $\forall i \in \{1, \cdots, s\}$, $z_i \in \mathbb{R}^d$ is replaced by its nearest neighbor item $m_{\mathrm{J}} \in \mathbb{R}^d$ in the memory as follows:

$$\tilde{z}_i \leftarrow m_{\mathrm{J}}, \ \text{where J} = \text{argmin}_j \|z_i - m_j\|_2, \quad j \in \{1, \cdots, k\}. \tag{2}$$

Then the input to the decoder $\boldsymbol{Dec}_p$ is the substituted feature $\tilde{z}$. To simplify the description, we denote $\mathcal{G}$ as the function of retrieving the corresponding item in the memory and using the retrieved item to replace the input. Then we arrive at the following mapping function:
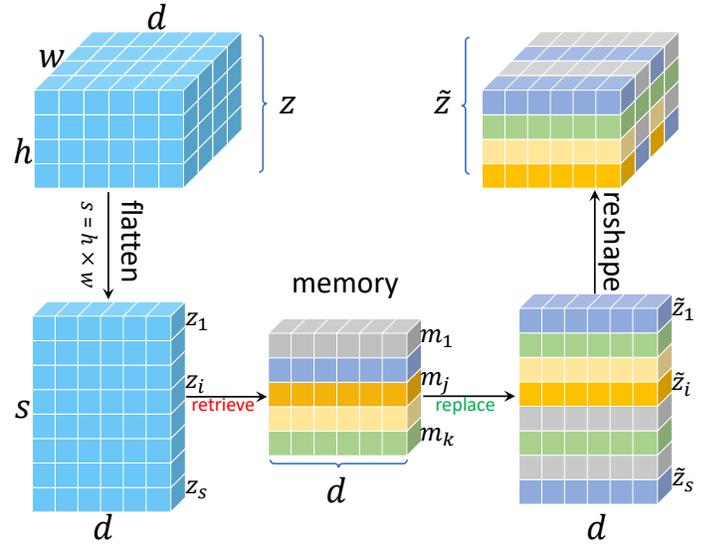
$$\tilde{z} = \mathcal{G}(z) \tag{3}$$



Fig. 3. The illustration of the detailed operation in memory. $z_i \in \mathbb{R}^d$ retrieves the nearest vector $m_{\mathrm{J}} \in \mathbb{R}^d$ in the memory, and $z_i$ is replaced by $\tilde{z}_i$ with Equation (2).

In the training phase, we update the memory with exponential moving averages [49]. We denote $z_{i,1}, z_{i,1}, \ldots, z_{i,n_i}$ as $n_i$ latent features that are closest to the memory item $m_i$. In order to make the memory item close to the set of latent features, we have the loss as:

$$\sum_{j}^{n_i} \|z_{i,j} - m_i\|^2. \tag{4}$$

The optimal $m_i$ has a closed form solution, which is the average of the latent features in the set:

$$m_i = \frac{1}{n_i} \sum_{j}^{n_i} z_{i,j}. \tag{5}$$

In order to handle online mini-batch training, the exponential moving average [49] is used:

$$N_i^{(t)} = N_i^{(t-1)}\gamma + n_i^{(t)}(1 - \gamma)$$
$$e_i^{(t)} = e_i^{(t-1)}\gamma + \sum_j z_{i,j}^{(t)}(1 - \gamma)$$
$$m_i^{(t)} = \frac{e_i^{(t)}}{N_i^{(t)}},$$
(6)

where $\gamma$ is a constant between 0 and 1, meaning how much history data to be kept. $n_i^{(t)}$ and $z_{i,j}^{(t)}$ denote $n_i$ and $z_{i,j}$ in the $t-th$ mini-batch, respectively.

During the training of Proxy Extraction Module, we update the encoder and decoder simultaneously with the update of the memory. The encoder and decoder are optimized as:

$$\mathcal{L}_p = \|\boldsymbol{Dec}_p(\mathcal{G}(\boldsymbol{Enc}_p(\mathbf{I}))) - \mathbf{P}\|_2^2,$$
(7)

where $\mathbf{P}$ is the SI extracted from the original image $\mathbf{I}$ with the SLIC algorithm [48]. It is worth noting that there is no gradient defined for the function $\mathcal{G}$. However, it is possible to approximate the gradient with the straight-through estimator [49]. We copy the gradient of $\tilde{z}$ as the gradient of $z$. Once the Proxy Extraction Module is trained, the weights of $\boldsymbol{Dec}_p$, $\boldsymbol{Enc}_p$ and the items in memory are fixed when optimizing the Image Reconstruction Module.

**Our method vs. MemAE [50]**. Recently, Gong *et al.* [50] proposed to augment an Auto-Encoder with the memory module, termed MemAE. Given an input, the encoder in MemAE extracts an encoded representation, which is used as a query to retrieve the most relevant items in the memory. The multiple items are then averaged with an attention weight to get the substituted feature $\tilde{z}$ for the reconstruction. Both the MemAE and our method can take advantage of memory-augmented networks for anomaly detection. However, the specific task in this work is relatively simpler than [50]. Thus, our method only retrieve the nearest item as the substituted feature rather that retrieve multiple items. The reasons why the task in this paper is relatively simpler are in two aspects: first, the patterns in the medical images are simpler than the patterns in the natural images [51], [52], [53]; second, MemAE [50] reconstructs original image while our method predicts the SI, which contains the simple pattern than that in original image. As a result, if a combination of multiple memory items is used, some anomalies may still have the chance to be well reconstructed. Therefore, we only retrieve a single item rather multiple items for SI prediction.

### D. Image Reconstruction Module

In anomaly detection, the abnormal images are usually not available in the training phase. However, compared with the normal SI, the input image could be regarded as a pseudo abnormal SI in appearance. This inspires us to synthesize a pseudo abnormal proxy by cutting a patch from the normal image, and paste it into the normal proxy. Then we obtain the pseudo abnormal proxy $\mathbf{P}'$, and $\mathbf{M} \in \{0, 1\}$ denotes a binary mask indicating the location of pasting the patch into $\mathbf{P}'$. The pseudo abnormal proxy is inputted into the Image

Reconstruction Module and we propose a repairing loss to enforce the output of abnormal proxy as a normal image. Thus, in the test phase, the reconstruction error corresponding to the abnormal images would be boosted, and consequently, the discrepancy between the normal images and the abnormal images is enlarged, which is a desirable feature for anomaly detection.

To obtain diverse pseudo abnormal proxy $\mathbf{P}'$, we propose a Pseudo Abnormal Proxy Constructor (PAPC) that crops the patch from the image with a random size and pasts it into the proxy at a random position. As shown in Fig. 4, the pseudo abnormal proxy can mimic the abnormal images to some extent. In the training phase, we enforce the network to reconstruct the normal image even with the pseudo abnormal SI, which is like the 'anomalies repairing'. In this way, the reconstruction error corresponding to the abnormal images would be boosted, and consequently, the discrepancy between the normal images and the abnormal images is enlarged, which is a desirable feature for anomaly detection.
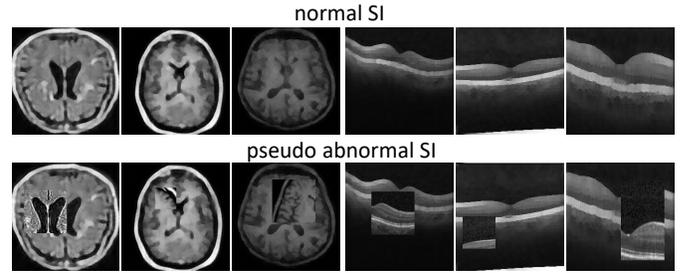


Fig. 4. Examples of the normal SI (first row) and the pseudo abnormal SI (second row) in brain MRI and retinal OCT.

Specifically, in the training of the Image Reconstruction Module, the forward pass of the normal proxy and the pseudo abnormal proxy is denoted as $\hat{\mathbf{I}} = \mathcal{F}_g(\hat{\mathbf{P}})$ and $\hat{\mathbf{I}}' = \mathcal{F}_g(\mathbf{P}')$, respectively. The function $\mathcal{F}_g(\cdot)$ is shared.

**Normal Image Reconstruction Loss.** For the normal proxy, the image reconstruction loss is shown as follow:

$$\mathcal{L}_{\text{rec}} = \|\hat{\mathbf{I}} - \mathbf{I}\|_2^2 + \lambda_g \underbrace{\big(\mathbb{E}[\log(1 - \mathbf{D}(\hat{\mathbf{I}}))] + \mathbb{E}[\log\mathbf{D}(\mathbf{I})]\big)}_{\text{adversarial loss}},$$
(8)

where $\mathbf{D}$ denotes the discriminator [16], and the $\lambda_g = 0.01$ is a hyper-parameter.

**Abnormal Image Repairing Loss.** For the pseudo abnormal proxy, we propose a repairing loss to enforce the output of abnormal proxy as a normal image. The repairing loss consists of a global regularization item and a local regularization item, which are defined as follows:

$$\mathcal{L}_{\text{global}} = \|\hat{\mathbf{I}}' - \mathbf{I}\|_2^2 + \lambda_g\big(\mathbb{E}[\log(1 - \mathbf{D}(\hat{\mathbf{I}}'))] + \mathbb{E}[\log\mathbf{D}(\mathbf{I})]\big),$$
(9)

$$\mathcal{L}_{\text{local}} = \|\mathbf{M} \odot \hat{\mathbf{I}}' - \mathbf{M} \odot \mathbf{I}\|_2^2 + \lambda_g\big(\mathbb{E}[\log(1 - \mathbf{D}(\mathbf{M} \odot \hat{\mathbf{I}}'))] + \mathbb{E}[\log\mathbf{D}(\mathbf{M} \odot \mathbf{I})]\big),$$
(10)

where $\odot$ is element-wise multiplication.

**Total Objective Function.** We arrive at the objective function for the Image Reconstruction Module:

$$\begin{aligned}\mathcal{L} &= \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{repairing}} \\ &= \mathcal{L}_{\text{rec}} + \lambda_{\text{global}}\mathcal{L}_{\text{global}} + \lambda_{\text{local}}\mathcal{L}_{\text{local}},\end{aligned} \quad (11)$$

where $\lambda_{\text{gobal}}$ and $\lambda_{\text{local}}$ are the hyper-parameters. Empirically, we set $\lambda_{\text{gobal}} = 0.25$, $\lambda_{\text{local}} = 0.5$ on all datasets in our experiments.

### E. Anomaly Detection on Test Data

Taking an original image $\mathbf{I} \in \mathbb{R}^{w \times h}$ as an input, the proposed method can obtain an extracted proxy $\hat{\mathbf{P}} \in \mathbb{R}^{w \times h}$ and a reconstructed image $\hat{\mathbf{I}} \in \mathbb{R}^{w \times h}$. The $w$ and $h$ denote the width and the height of image, respectively. We compute the anomaly score for the image-level anomaly detection in the latent feature space, and compute the anomaly score map for the pixel-level anomaly detection in the image space.

**Anomaly Score for Image-level Anomaly Detection.** It has been proven that computing the anomaly score by mapping the image space to the latent space is effective [18], [36]. However, training an additional encoder is inefficient and redundant. We apply the existing encoder $\textit{\textbf{Enc}}_p$ in Proxy Extraction Module to map the image to a latent space, and we get the latent feature as:

$$\begin{aligned} z &= \textit{\textbf{Enc}}_p(\mathbf{I}) \\ \hat{z} &= \textit{\textbf{Enc}}_p(\hat{\mathbf{I}}) \end{aligned}, \quad (12)$$

where $\hat{z}$ is the latent feature of the reconstructed image $\hat{\mathbf{I}}$.

We compute the anomaly score as:

$$\mathbf{A}_{\text{img}} = \|z - \hat{z}\|_F. \quad (13)$$

**Anomaly Score Map for Pixel-level Anomaly Detection.** To get the pixel-level anomaly map $\mathbf{A}_{\text{pix}} \in \mathbb{R}^{w \times h}$ for lesion segmentation, we compute the anomaly score map in the image space as:

$$\mathbf{A}_{\text{pix}} = |\mathbf{I} - \hat{\mathbf{I}}|. \quad (14)$$

### F. Detailed Network Architecture

We use the same encoder and decoder in both the Proxy Extraction Module and the Image Reconstruction Module. The detailed architectures of the encoder and the decoder are shown in Fig. 5.
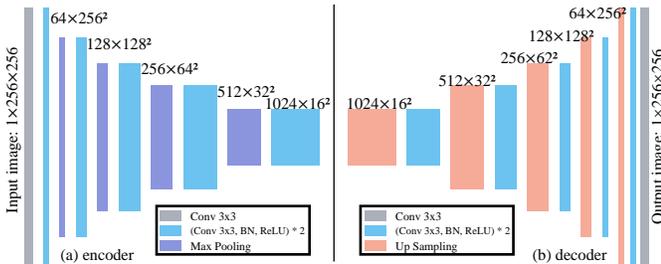


Fig. 5. The architectures of the encoder and decoder.

## IV. EXPERIMENTS

### A. Experimental Setup

*1) Implementation Details:* The optimizer we used is the Adam optimizer [57], and the learning rate is set as 0.001. The scale of original image and its proxy are resized to be $256 \times 256$. The number of superpixels is set to be 800 in the SLIC [48] algorithm. In the memory, we set $d = 64$, $k = 128$, $\gamma = 0.99$. In our implementation, the Proxy Extraction Module and the Image Reconstruction Module are trained separately, and we find that training two modules within an end-to-end learning manner does not necessarily bring performance gain.

*2) Evaluation Metrics:* For performance evaluation, we calculate the Area Under Receiver Operation Characteristic (AUC) by gradually changing the threshold of anomaly score $\mathbf{A}_{\text{pix}}$ and $\mathbf{A}_{\text{img}}$ for the pixel-level and the image-level anomaly detection, respectively. By following [36], [58], we also calculate the accuracy (ACC) and F1-score for the performance evaluation. The F1-score is defined as the harmonic mean of precision and recall of a model.

*3) Baseline Methods:* We compare our method with several state-of-the-art anomaly detection methods, including MemAE [1] [50], AnoGAN [2] [4], f-AnoGAN [3] [35], GANomaly [4] [18], Sparse-GAN [36], Auto-Encoder [14] based anomaly detection method, and pix2pix [5] [43] based anomaly detection method. As our method consists of the translation between image and proxy, we also compare our ProxyAno with Cycle-GAN [5] [56]. We use the code provided in the GitHub to implement the baseline methods. For a fair comparison, the architecture of the encoder and decoder in Auto-Encoder is the same as that used in our method.

### B. Image-level Anomaly Detection on Brain MRI

*1) **Dataset**:* The training dataset used in previous brain MRI anomaly detection work [14], [17] are not released. However, there are several publicly available brain MRI datasets [28], [59] used for tumor classification and lesion segmentation. Thus, we propose to integrate [28] and [59] and use it to evaluate our proposed method. There are 500 normal images in the training set, and the test set consists of 500 abnormal image and 500 normal images.

*2) **Performance Evaluation**:* We report the performance of different methods for image-wise anomaly detection on the brain MRI dataset in Table I. We can see that the AUC and ACC of our proposed ProxyAno outperforms that of all baseline methods on the brain MRI dataset. In particular, ProxyAno achieves an AUC of 0.853 on the brain MRI dataset with an increase of 1.8% from 0.835 by the latest Sparse-GAN [36]. However, the F1-score of our method does not achieve the best. From Table I, the similar results can be found in other two datasets [54][55], *i.e.*, the AUC of our method achieves the best while the F1-score of our method does not. The reason

---

[1] https://github.com/donggong1/memae-anomaly-detection
[2] https://github.com/LeeDoYup/AnoGAN-tf
[3] https://github.com/tSchlegl/f-AnoGAN
[4] https://github.com/samet-akcay/ganomaly
[5] https://github.com/junyanz/pytorch-CycleGAN-and-pix2pix

TABLE I
The image-level anomaly detection results on brain MRI and retinal OCT, and the pixel-level anomaly detection results on retinal fundus. The best two results are shown in **RED** and RED fonts.

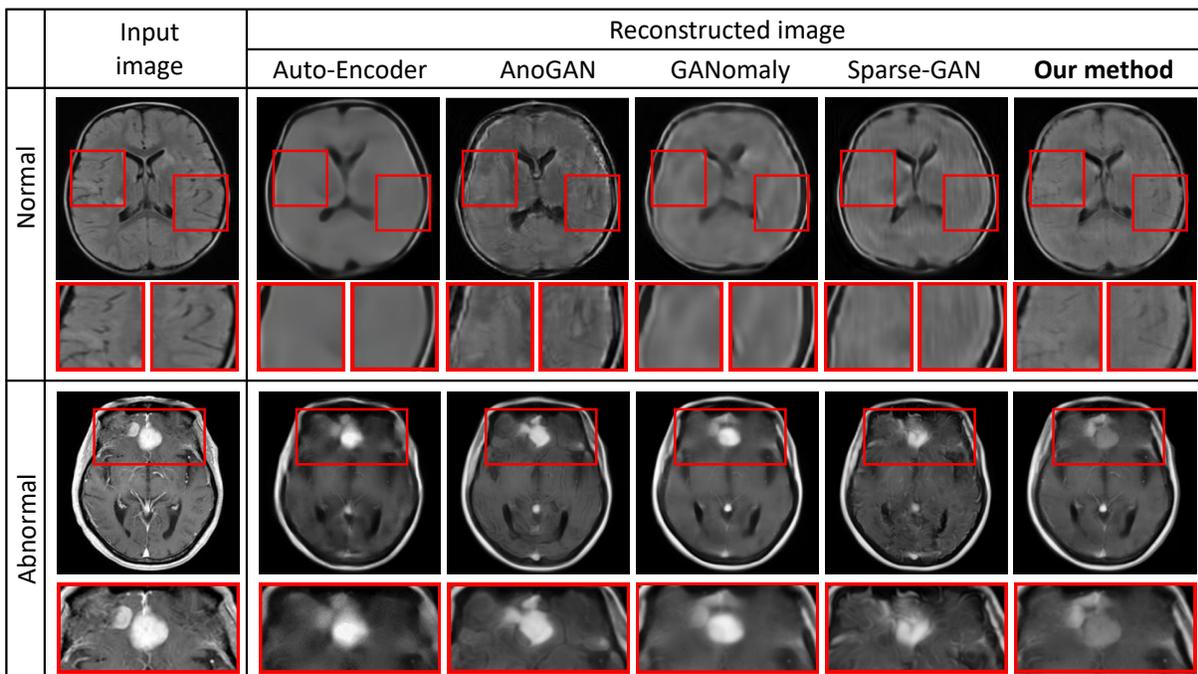| Mehtods | Brain MRI [28] | | | Retinal OCT [54] | | | Retinal Fundus [55] | | |
|---|---|---|---|---|---|---|---|---|---|
| | AUC | ACC | F1-score | AUC | ACC | F1-score | AUC | ACC | F1-score |
| Auto-Encoder [14] | 0.766 | 0.714 | 0.674 | 0.783 | 0.751 | 0.669 | 0.609 | 0.589 | 0.542 |
| MemAE [50] | 0.848 | 0.789 | **0.722** | 0.876 | 0.838 | 0.652 | 0.647 | 0.592 | 0.567 |
| AnoGAN [4] | 0.803 | 0.757 | 0.691 | 0.846 | 0.789 | 0.637 | 0.630 | 0.618 | 0.579 |
| f-AnoGAN [35] | 0.822 | 0.764 | 0.675 | 0.882 | 0.808 | 0.653 | 0.698 | **0.686** | 0.637 |
| GANomaly [18] | 0.832 | 0.798 | 0.667 | 0.916 | 0.826 | **0.727** | 0.652 | 0.633 | **0.658** |
| Sparse-GAN [36] | 0.835 | 0.791 | 0.645 | 0.925 | 0.841 | 0.714 | 0.663 | 0.638 | 0.651 |
| pix2pix [43] | 0.796 | 0.737 | 0.617 | 0.861 | 0.818 | 0.702 | 0.632 | 0.621 | 0.603 |
| Cycle-GAN [56] | 0.808 | 0.752 | 0.712 | 0.815 | 0.762 | 0.675 | 0.626 | 0.613 | 0.608 |
| **Our Method** | **0.853** | **0.805** | 0.709 | **0.933** | **0.849** | 0.725 | **0.701** | 0.682 | 0.649 |



Fig. 6. The reconstruction results on the normal and the abnormal brain MRI. The red rectangle in the abnormal sample denotes the brain tumor. From left to right: input brain MRI, reconstruction result obtained by Auto-Encoder [14], AnoGAN [4], GANomaly [18], Sparse-GAN [36] and our method. Our method reconstructs the abnormal regions poorly while reconstructing the normal image well.

behind this is probably the threshold selection has an important effect on the F1-score while it dose not effect the AUC result.

We also show the reconstruction results in Fig. 6 to visually compare our method with other competitive methods, including Auto-Encoder [14], AnoGAN [4], GANomaly [18] and Sparse-GAN [36]. The images show that our method can well reconstruct the normal images while poorly reconstruct the abnormal ones.

### C. Image-level Anomaly Detection on Retinal OCT

*1) Dataset:* The retinal OCT [54] dataset from Spectralis OCT (Heidelberg Engineering, German) is used in this section, and it contains data with three different lesions: drusen, DME (diabetic macular edema), and CNV (choroidal neovascularization). This dataset contains the standard training/test split. We use the normal images in the original training set to train the model, and use all the test images for performance evaluation.

*2) Performance Evaluation:* As shown in the Table I, our method outperforms all baseline methods on the retinal OCT

dataset. Particularly, the proposed ProxyAno achieves an AUC of 0.933 in the retinal OCT. Further, our method significantly outperforms the Auto-Encoder based solution because our solution can mitigate the identity mapping issue and increase the reconstruction error for the abnormal images. Our method also outperforms the Cycle-GAN based solution because we introduce the pseudo abnormal SI's, which would be repaired as the normal images. Consequently, the reconstruction errors for the abnormal images would be large, which would also facilitate the anomaly detection.

We further compute the average anomaly score for the image-level anomaly detection for both the normal images and the abnormal images on the OCT dataset with Equation (13). Then, we calculate the gap of these two scores to measure the ability of our method and other baselines to discriminate the normal and the abnormal images. A larger gap means the normal and the abnormal images can be more easily separated. The results in Table II show that our method achieves the largest gap than other baselines, which validates

the effectiveness of our method for anomaly detection.

TABLE II
THE AVERAGE ANOMALY SCORE FOR THE NORMAL IMAGES, THE
ABNORMAL IMAGES AND THE GAP BETWEEN THESE TWO SCORES ON THE
RETINAL OCT DATASET. OUR SI APPROACH DENOTES USING SI TO
CONNECT TWO PRIMITIVE ENCODER-DECODER PAIRS, WITHOUT USING
MEMORY AND REPAIRING LOSS.

| Method | Anomaly Score | | |
|---|---|---|---|
| | normal | abnormal | gap |
| Auto-Encoder [14] | 0.601 | 0.898 | 0.297 |
| f-AnoGAN [35] | 0.549 | 0.922 | 0.373 |
| GANomaly [18] | 0.430 | 0.834 | 0.404 |
| Sparse-GAN [36] | 0.576 | 0.995 | 0.419 |
| Our SI approach | 0.580 | 1.006 | 0.426 |
| Our Final Model | 0.568 | 1.003 | **0.435** |

### D. Pixel-level Anomaly Detection on Retinal Fundus Image

*1) Dataset*: Indian Diabetic Retinopathy Image Dataset (IDRiD) [55] is a publically available dataset used to evaluate Diabetic Retinopathy (DR) detection. To train the proposed model, we only choose the normal class from the original training set in IDRiD as our training set, and the test set is the same as the original lesion detection dataset. IDRiD consists of 134 normal images for training, and 69 abnormal images with pixel-level lesion annotation for testing. Since the resolution of original image is very large (i.e., $4288 \times 2848$), we crop 9 (i.e., $3 \times 3$) patches from each image.

*2) Performance Evaluation*: The baseline methods for pixel-level anomaly detection are the same as that in image-level performance evaluation. To compute the lesion detection accuracy, we set the threshold of $\mathbf{A}_{pix}$ as 0.5. As shown in Table I, the AUC and ACC of proposed method outperforms the baseline methods, which further verifies the effectiveness of our method. We further qualitatively show the reconstruction and lesion detection results in Fig. 7. It can be found that the normal region is reconstructed with a small error, while the region with lesion is reconstructed with a large error. Both the red lesions and bright lesions can be roughly detected. Here the red lesions include haemorrhages of all shapes and microaneurysms, bright lesions include hard and soft exudates, drusen, cotton-wool spots [60].

### E. The Selection of The Proxy

In this part, to verify the effectiveness of proxy approach, we study different proxy types and compare them with baseline that without proxy (i.e., Auto-Encoder). Specifically, besides the SI, we also consider the following five types of proxy: 1) the edge extracted by a Canny edge detector; 2) smooth image, which is a short for the image smoothed with Gaussian blur; 3) image with smooth patches, which denotes the image consists of patches colored with the average intensity in each patches; 4) edge $\oplus$ smooth image, here $\oplus$ denotes the concatenation; 5) edge $\oplus$ image with smooth patches. The experiments are conducted on the retinal OCT [54] dataset. The illustration of smooth image and image with smooth patches can be found in Fig. 8.

The quantitative results are reported in Table III, from which we can see that all proxy methods outperform the baseline

without proxy, validating proxy approach is effective. We can also see that the SI-based proxy is superior to all other types of proxy. The reason of SI-based proxy outperforms the edge-based proxy is possibly because that the edge contains less texture information, and makes the image reconstruction difficult. Compared with the edges, the proposed SI contains more texture information (i.e., the average intensity within each superpixel), which makes the image reconstruction easier, and facilitates the anomaly detection. The reason of SI-based proxy outperforms the proxy based on smooth image and image with smooth patches is possibly because that the SI contains more structure information. Compare with the smooth image, in the SI the edge is relatively enhanced. Compare with the image with smooth patches, the SI contains more semantic information and preserves local structure. As a summary, compare with other types of proxy, the proposed SI contains more texture and structure information, validating that both texture and structure information are important for anomaly detection. This similar argument can also be found in [47]. Moreover, with concatenating edge, the performance of both smooth image and image with smooth patches improved. This also validate the argument that both texture and structure information are important for anomaly detection.

TABLE III
THE RESULTS OF DIFFERENT PROXY SELECTION STRATEGIES.

| Index | Different Proxies | AUC |
|---|---|---|
| 0 | without proxy (Auto-Encoder [14]) | 0.783 |
| 1 | edge | 0.804 |
| 2 | smooth image | 0.801 |
| 3 | image with smooth patches | 0.796 |
| 4 | edge $\oplus$ smooth image | 0.813 |
| 5 | edge $\oplus$ image with smooth patches | 0.815 |
| 6 | **proposed SI** | **0.818** |

The qualitative results are shown in Fig. 9. As shown in the images in the $1^{st}$ row of Fig. 9, compared with the edge-based reconstruction, our SI-based solution reconstruct the normal image more accurately. Furthermore, the images in the $2^{nd}$ column of Fig. 9 show that although the self-reconstruction based method (i.e., 'image $\rightarrow$ image' using the AE) is trained well on the normal OCT images, when we feed the abnormal images (i.e., DME image and drusen image) to the trained AE, it can also well reconstruct the lesions in the abnormal images. On the contrast, if we use an SI as the proxy, the issue is eased. Specifically, as shown by the images in the $5^{th}$ column, since the SI is obtained by a neural network that trained on the normal samples, the network cannot well extract the SI on the DEM images and the drusen images, especial the lesion regions in the abnormal images. Based on the poorly extracted SI, the abnormal images cannot be well reconstructed. The images in the $6^{th}$ column show that the lesions in the DEM image and the drusen image tend to be repaired, thus the reconstruction error between input image and reconstructed image is larger than self-reconstruction based method.

### F. Ablation Studies

In this part, we conduct several ablation studies on retinal OCT dataset [54] to evaluate each component of our method.
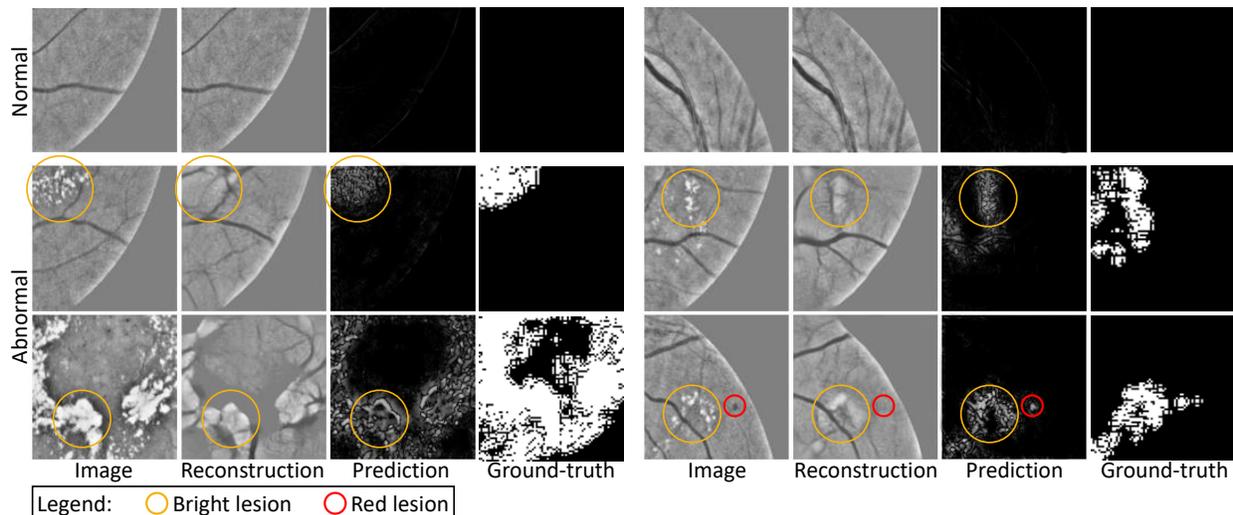
Fig. 7. The qualitative results on the retinal fundus image. Both the red lesions and bright lesions can be roughly detected with our method.
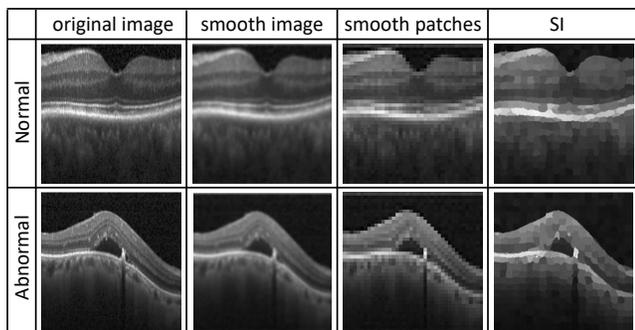


Fig. 8. Illustration of different types of proxy. Compared with smooth image and smooth patches (short for image with smooth patches), the SI preserves local structure.

TABLE IV
THE ABLATION STUDIES OF PROPOSED METHOD.

| Index | Method | AUC |
|---|---|---|
| 1 | EncDec (Auto-Encoder [14]) | 0.783 |
| 2 | EncDec + mem* (MemAE [50]) | 0.816 |
| 3 | EncDec + mem | 0.820 |
| 4 | 2×EncDec + SI | 0.818 |
| 5 | 2×EncDec + SI + mem | 0.854 |
| 6 | 2×EncDec + SI + rep | 0.869 |
| 7 | 2×EncDec + SI + mem + rep | 0.916 |
| 8 | 2×EncDec + SI + mem + rep + lat (final) | **0.933** |

We report the results of the difference components in Table IV. Let EncDec denotes the encoder-decoder architecture. The relationship between EncDen and Auto-Encoder is: if the output of EncDec is the same as the input, the EncDec is actually the Auto-Encoder; otherwise, the EncDec is not the Auto-Encoder. The mem* denotes the memory implemented in [50], while the mem denotes the memory implemented in this work. The specific difference between mem* and mem is described in the Section III-C, Proxy Extraction Module. "EncDec + mem" (or mem*) denotes that there is a memory between the encoder and decoder in the Auto-Encoder. Both Proxy Extraction Module and Image Reconstruction Module contain the EncDec, thus we use "2×EncDec" to simply denote Proxy Extraction Module and Image Reconstruction Module. Between these two modules, we use SI as the intermediate proxy, which is denoted as "2×EncDec + SI". Based on "2×EncDec + SI", "+ mem" denotes the Proxy Extraction Module is equipped with the memory, while "+ rep" denotes the Image Reconstruction Module is equipped with the proposed repairing loss. "+ lat" denotes that we compute the anomaly score in the latent space rather than in the image space.

*1) The effectiveness of memory in Auto-Encoder:* We first validate the effectiveness of the memory between the encoder and decoder in Auto-Encoder. As shown in row 1, row 2, and row 3 in Table IV, both memory implemented by [50] and implemented by ours outperform the Auto-Encoder, validating that explicitly memorize the correspondence in the feature space with memory is effective. Besides, the results of the memory implemented by ours is better than that implemented by [50], validating that our memory is more suitable for anomaly detection in medical image. The reason behind this is: MemAE [50] retrieves multiple items while our memory only retrieves the nearest item, and combining multiple items result in some anomalies may still being well reconstructed; therefore, our memory is better than MemAE [50] for anomaly detection in medical image.

*2) The effectiveness of using SI as the proxy:* The results in Table IV show that using SI as the intermediate proxy (both with and without memory) consistently achieves higher AUC than Auto-Encoder. By comparing row 4 versus row 1 and row 5 versus row 3, it can be found that using SI as the intermediate proxy can roughly improve 3.5% AUC. This proves that it is necessary to use SI as the proxy. Additionally, the results in row 4 and row 2 show that our SI approach outperforms MemAE [50].

*3) The effectiveness of memory in Proxy Extraction Module:* Based on "2×EncDec + SI", we equip the Proxy Extraction Module with memory to verify its effective. As shown the row 5 in Table IV, adding the memory increases the performance. Furthermore, we show the qualitative results
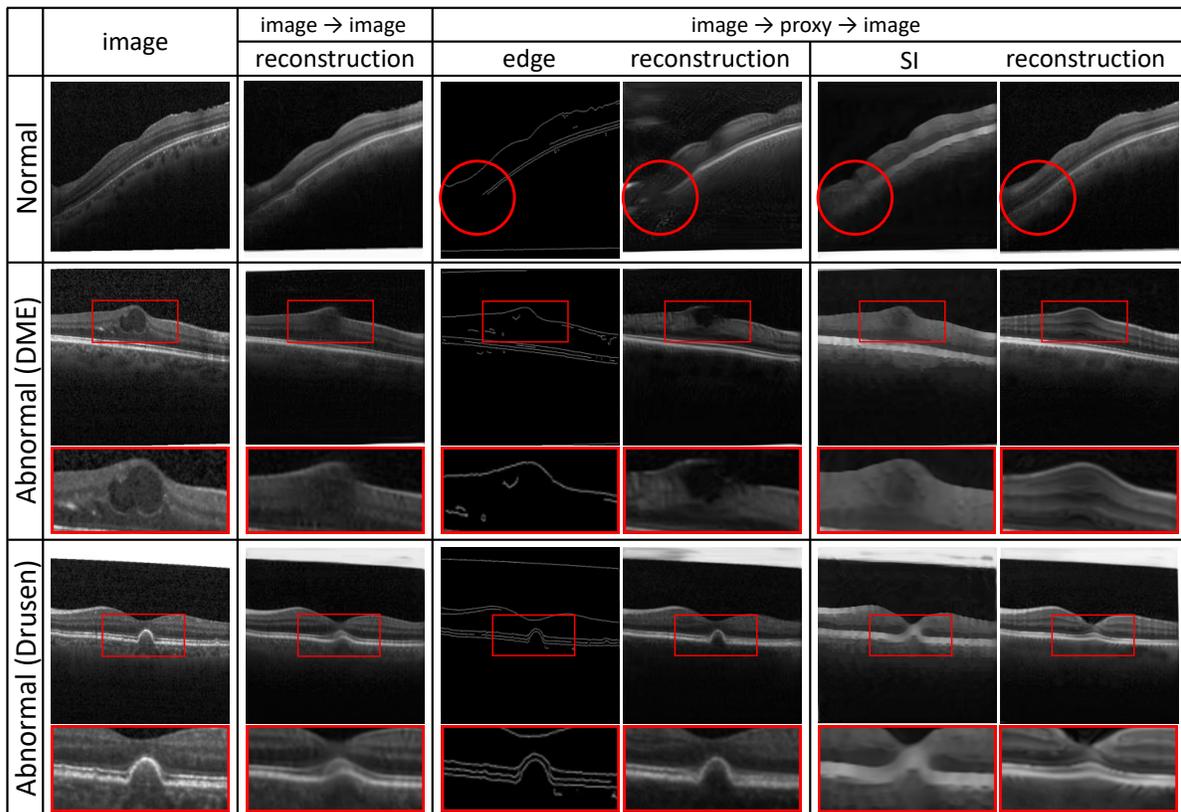
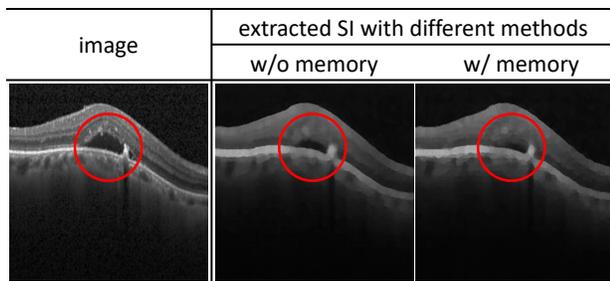Fig. 9. The visualization results with different proxy.



Fig. 10. The qualitative results of extracted SI with different methods. The red circle denotes the DME lesion. Our method tries to fill the lesion hole with normal patches. Consequently, our solution can easily identify the abnormal images.

TABLE V
THE AUC RESULTS OF DIFFERENT WAYS FOR ANOMALY SCORE
CALCULATING.

|  | Brain MRI | Retinal OCT |
|---|---|---|
| image space | 0.807 | 0.916 |
| latent feature space | **0.853** | **0.933** |

of extracted SI in Fig. 10, and it can be observed that the DME lesion in the memory-based method is similar to the normal SI. This agrees with our assumption that memory makes the abnormal images look like the normal ones, which enlarges the reconstruction error of the abnormal images and facilitates the anomaly detection.

*4) **The effectiveness of the repairing loss in Image Reconstruction Module**:* The results in Table IV show that using the proposed repairing loss (both with and without memory) consistently achieves higher AUC than that without repairing loss. By comparing row 6 versus row 4 and row 7 versus row 5, it can be found that trained with repairing loss can roughly improve 5% AUC. This proves the effectiveness of the repairing loss in Image Reconstruction Module.

*5) **The way of calculating anomaly score for image-level anomaly detection**:* In this part, we compare the way of calculating anomaly score in image space with that in latent feature space. Given an input image $\mathbf{I}$ and its reconstructed counterpart $\hat{\mathbf{I}}$, we can calculate the anomaly score in the **image space** based on the difference between them ($\|\mathbf{I} - \hat{\mathbf{I}}\|_F$). We compute the anomaly score in the **latent feature space** with Equation (13). The results are reported in Table V, we can see that the anomaly score calculated in the latent feature space always corresponds to the better performance. The possible reason is that the reconstructed image always contain some noises, which may affect the anomaly detection. While the latent feature can be less affected by the noises, and boosts the anomaly detection. It is worth noting that previous work [18], [35], [36] also use the latent feature to calculate the anomaly score.

### G. Hyper-parameters Analysis

*1) **The effect of memory size**:* We gradually change the memory size and show the results in Fig. 11, and it shows that when $k = 128$ the model achieves best performance. The
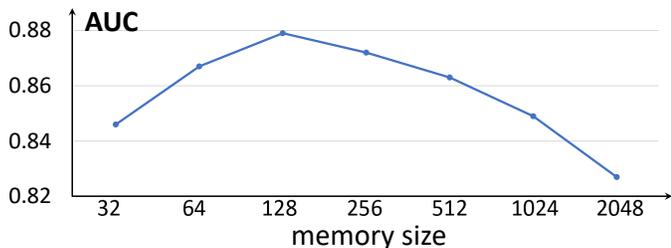
Fig. 11. The AUC results of our Proxy Extraction Module with different memory size. When setting the memory size as 2048, the model degenerates to an Encoder-Decoder model without memory.

memory with a small size is incapable to characterize the mapping between the normal input and its SI, while the memory with large size cannot repair the tissues in abnormal images in the SI extraction. Theoretically, a memory with an infinite size corresponds to a model without the memory. Particularly, when setting the memory size as 2048, the performance of AUC decreased to 0.837, which is close to the baseline model (Encoder-Decoder w/o memory) that achieves 0.828. This also verifies the effectiveness of memory for anomaly detection.
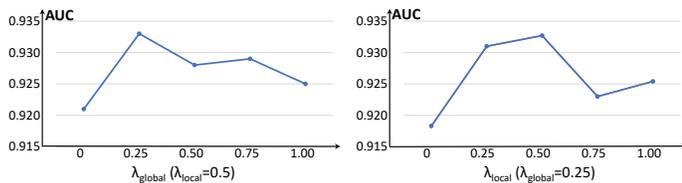


Fig. 12. The AUC results of the model with different weights.

*2) The pseudo anomaly repairing loss:* In our implementation, the weight for the global repairing term and local repairing term is 0.25 and 0.5, respectively. Then we fix the weight of one term and change the weight of another term. As shown in the in Fig. 12, we can see that the results with only local term ($\lambda_{\text{global}} = 0$) or global term ($\lambda_{\text{local}} = 0$) are worse than the results based on both terms. We further show the qualitative results in Fig. 13, and it can be observed that the lesion is repaired when both global and local terms are used, which verifies the effectiveness of these terms for the anomaly detection.
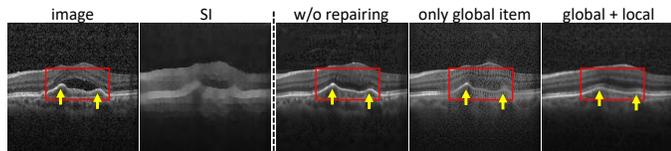


Fig. 13. The visualization results of the model with different training loss. The red rectangle denotes the DME lesion.

## V. CONCLUSION

In this paper, we propose a novel ProxyAno approach to mitigate the identity mapping issue in the Auto-Encoder based anomaly detection paradigm. Specifically, our ProxyAno first maps an input image to a superpixel-images based on a memory-aided Proxy Extraction Module, and this memory would cause information loss for abnormal images and facilitate the anomaly detection. Further, an Image Reconstruction Module is used to reconstruct the input based on the

superpixel-images. To further enlarge the reconstruction error, we enforce the reconstruction module can well reconstruct the input even with a pseudo abnormal superpixel-image. In this way, our ProxyAno favors the reconstruction of normal input and leads to a large reconstruction of the abnormal images. Extensive experiments on different medical image datasets validate the effectiveness of our approach for both image-level and pixel-level anomaly detection.

## VI. ACKNOWLEDGE

## REFERENCES

[1] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly detection: A survey," *ACM Computing Surveys (CSUR)*, vol. 41, no. 3, pp. 1–58, 2009.

[2] G. Pang, C. Shen, L. Cao, and A. V. D. Hengel, "Deep learning for anomaly detection: A review," *ACM Computing Surveys (CSUR)*, vol. 54, no. 2, pp. 1–38, 2021.

[3] 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 Transactions on Neural Networks and Learning Systems*, 2021.

[4] 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 *International Conference on Information Processing in Medical Imaging*. Springer, 2017, pp. 146–157.

[5] D. Zimmerer, F. Isensee, J. Petersen, S. Kohl, and K. Maier-Hein, "Unsupervised anomaly localization using variational auto-encoders," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2019, pp. 289–297.

[6] C. Baur, B. Wiestler, S. Albarqouni, and N. Navab, "Bayesian skip-autoencoders for unsupervised hyperintense anomaly detection in high resolution brain mri," in *2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI)*. IEEE, 2020, pp. 1905–1909.

[7] X. Chen, S. You, K. C. Tezcan, and E. Konukoglu, "Unsupervised lesion detection via image restoration with a normative prior," *Medical image analysis*, p. 101713, 2020.

[8] X. Li, X. Hu, L. Yu, L. Zhu, C.-W. Fu, and P.-A. Heng, "Canet: Cross-disease attention network for joint diabetic retinopathy and diabetic macular edema grading," *IEEE transactions on medical imaging*, vol. 39, no. 5, pp. 1483–1493, 2019.

[9] L. Luo, L. Yu, H. Chen, Q. Liu, J. Xu, P.-A. Heng *et al.*, "Deep mining external imperfect data for chest x-ray disease screening," *IEEE transactions on medical imaging*, 2020.

[10] P. P. Srinivasan, L. A. Kim, P. S. Mettu, S. W. Cousins, G. M. Comer, J. A. Izatt, and S. Farsiu, "Fully automated detection of diabetic macular edema and dry age-related macular degeneration from optical coherence tomography images," *Biomedical optics express*, vol. 5, no. 10, pp. 3568–3577, 2014.

[11] K. Zhou, Z. Gu, W. Liu, W. Luo, J. Cheng, S. Gao, and J. Liu, "Multi-cell multi-task convolutional neural networks for diabetic retinopathy grading," in *2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*. IEEE, 2018, pp. 2724–2727.

[12] K. Ouardini, H. Yang, B. Unnikrishnan, M. Romain, C. Garcin, H. Zenati, J. P. Campbell, M. F. Chiang, J. Kalpathy-Cramer, V. Chandrasekhar *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.

[13] I. Golan and R. El-Yaniv, "Deep anomaly detection using geometric transformations," in *NeurIPS*, 2018.

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

[15] C. Zhou and R. C. Paffenroth, "Anomaly detection with robust deep autoencoders," in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2017, pp. 665–674.

[16] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.

[17] X. Chen and E. Konukoglu, "Unsupervised detection of lesions in brain mri using constrained adversarial auto-encoders," *Medical Imaging with Deep Learning (MIDL)*, 2018.

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

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

[20] H. Steck, "Autoencoders that don't overfit towards the identity," in *NeurIPS*, 2020.

[21] X. Ren and J. Malik, "Learning a classification model for segmentation," in *Computer Vision, IEEE International Conference on*, vol. 2. IEEE Computer Society, 2003, pp. 10–10.

[22] D. F. Kinane, "Causation and pathogenesis of periodontal disease," *Periodontology 2000*, vol. 25, no. 1, pp. 8–20, 2001.

[23] R. Nosalski and T. J. Guzik, "Perivascular adipose tissue inflammation in vascular disease," *British journal of pharmacology*, vol. 174, no. 20, pp. 3496–3513, 2017.

[24] L. Gutmann, D. Blumenthal, L. Gutmann, and S. S. Schochet, "Acute type ii myofiber atrophy in critical illness," *Neurology*, vol. 46, no. 3, pp. 819–821, 1996.

[25] 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 *31st British Machine Vision Conference*, 2020, pp. 7–10.

[26] D. Zimmerer, S. Kohl, J. Petersen, F. Isensee, and K. Maier-Hein, "Context-encoding variational autoencoder for unsupervised anomaly detection," in *International Conference on Medical Imaging with Deep Learning (MIDL)*, 2019.

[27] C. Baur, S. Denner, B. Wiestler, N. Navab, and S. Albarqouni, "Autoencoders for unsupervised anomaly segmentation in brain mr images: A comparative study," *Medical Image Analysis*, p. 101952, 2021.

[28] S. Rai, S. Chowdhury, S. Sarkar, K. Chowdhury, and K. P. Singh, "A hybrid approach to brain tumor detection from mri images using computer vision," *Journal of Innovation in Computer Science and Engineering*, vol. 8, no. 2, pp. 8–12, 2019.

[29] R. Bracewell, *Fourier analysis and imaging*. Springer Science & Business Media, 2004.

[30] D. A. Forsyth and J. Ponce, *Computer vision: a modern approach*. Prentice Hall Professional Technical Reference, 2002.

[31] G. Mori, X. Ren, A. A. Efros, and J. Malik, "Recovering human body configurations: Combining segmentation and recognition," in *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. IEEE, 2004.

[32] S. Kwak, S. Hong, and B. Han, "Weakly supervised semantic segmentation using superpixel pooling network," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 31, no. 1, 2017.

[33] D. Zimmerer, S. A. Kohl, J. Petersen, F. Isensee, and K. H. Maier-Hein, "Context-encoding variational autoencoder for unsupervised anomaly detection," *Medical Imaging for Deep Learning (MIDL)*, 2019.

[34] C. Han, L. Rundo, K. Murao, E. Sala, H. Nakayama, S. Satoh *et al.*, "Madgan: unsupervised medical anomaly detection gan using multiple adjacent brain mri slice reconstruction," *arXiv preprint arXiv:2007.13559*, 2020.

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

[36] K. Zhou, S. Gao, J. Cheng, Z. Gu, H. Fu, Z. Tu, J. Yang, Y. Zhao, and J. Liu, "Sparse-gan: Sparsity-constrained generative adversarial network for anomaly detection in retinal oct image," in *2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI)*. IEEE, 2020, pp. 1227–1231.

[37] P. Seebock, J. I. Orlando, T. Schlegl, S. M. Waldstein, H. Bogunovic, S. Klimscha, G. Langs, and U. Schmidt-Erfurth, "Exploiting epistemic uncertainty of anatomy segmentation for anomaly detection in retinal oct," *IEEE transactions on medical imaging*, 2019.

[38] J. Zhang, Y. Xie, G. Pang, Z. Liao, J. Verjans, W. Li, Z. Sun, J. He, Y. Li, C. Shen *et al.*, "Viral pneumonia screening on chest x-rays using confidence-aware anomaly detection," *IEEE transactions on medical imaging*, 2020.

[39] Y.-X. Tang, Y.-B. Tang, M. Han, J. Xiao, and R. M. Summers, "Abnormal chest x-ray identification with generative adversarial one-class classifier," in *2019 IEEE 16th International Symposium on Biomedical Imaging (ISBI 2019)*. IEEE, 2019, pp. 1358–1361.

[40] B. Bozorgtabar, D. Mahapatra, G. Vray, and J.-P. Thiran, "Salad: Self-supervised aggregation learning for anomaly detection on x-rays," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2020, pp. 468–478.

[41] N. Tuluptceva, B. Bakker, I. Fedulova, H. Schulz, and D. V. Dylov, "Anomaly detection with deep perceptual autoencoders," *arXiv preprint arXiv:2006.13265*, 2020.

[42] N. Pawlowski, M. C. Lee, M. Rajchl, S. McDonagh, E. Ferrante, K. Kamnitsas *et al.*, "Unsupervised lesion detection in brain ct using bayesian convolutional autoencoders," *Medical Imaging with Deep Learning (MIDL)*, 2018.

[43] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1125–1134.

[44] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 31, no. 1, 2017.

[45] 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," *International journal of computer vision*, vol. 115, no. 3, pp. 211–252, 2015.

[46] F. T. Liu, K. M. Ting, and Z.-H. Zhou, "Isolation forest," in *2008 eighth ieee international conference on data mining*. IEEE, 2008, pp. 413–422.

[47] 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 *ECCV*, 2020.

[48] R. Achanta, A. Shaji, K. Smith, A. Lucchi, P. Fua, and S. Süsstrunk, "Slic superpixels compared to state-of-the-art superpixel methods," *IEEE transactions on pattern analysis and machine intelligence*, vol. 34, no. 11, pp. 2274–2282, 2012.

[49] A. van den Oord, O. Vinyals, and K. Kavukcuoglu, "Neural discrete representation learning," in *Advances in Neural Information Processing Systems*, 2017, pp. 6306–6315.

[50] 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," *Proceedings of the IEEE International Conference on Computer Vision*, 2019.

[51] Y. Shen, B. Sheng, R. Fang, H. Li, L. Dai, S. Stolte, J. Qin, W. Jia, and D. Shen, "Domain-invariant interpretable fundus image quality assessment," *Medical image analysis*, vol. 61, p. 101654, 2020.

[52] Z. Shen, H. Fu, J. Shen, and L. Shao, "Modeling and enhancing low-quality retinal fundus images," *IEEE Transactions on Medical Imaging*, 2020.

[53] C.-H. Pham, A. Ducournau, R. Fablet, and F. Rousseau, "Brain mri super-resolution using deep 3d convolutional networks," in *2017 IEEE 14th International Symposium on Biomedical Imaging (ISBI 2017)*. IEEE, 2017, pp. 197–200.

[54] 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.

[55] P. Porwal, S. Pachade, M. Kokare, G. Deshmukh, X. Liu, L. Gao *et al.*, "Idrid: Diabetic retinopathy–segmentation and grading challenge," *Medical image analysis*, vol. 59, p. 101561, 2020.

[56] 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 IEEE international conference on computer vision*, 2017, pp. 2223–2232.

[57] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[58] J. Wolleb, R. Sandkühler, and P. C. Cattin, "Descargan: Disease-specific anomaly detection with weak supervision," *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2020.

[59] B. Sartaj, K. Ankita, B. Prajakta, and D. Sameer, *Brain Tumor Classification:Classify MRI images into four classes*, 2020. [Online]. Available: https://www.kaggle.com/sartajbhuvaji/brain-tumor-classification-mri

[60] C. Playout, R. Duval, and F. Cheriet, "A novel weakly supervised multitask architecture for retinal lesions segmentation on fundus images," *IEEE transactions on medical imaging*, vol. 38, no. 10, pp. 2434–2444, 2019.