

# In-painting Radiography Images for Unsupervised Anomaly Detection

Tiange Xiang<sup>1</sup> Yongyi Lu<sup>2</sup> Alan L. Yuille<sup>2</sup> Chaoyi Zhang<sup>1</sup> Weidong Cai<sup>1</sup> Zongwei Zhou<sup>2,\*</sup>  
<sup>1</sup>University of Sydney <sup>2</sup>Johns Hopkins University

## Abstract

We propose *space-aware memory queues* for *in-painting* and *detecting anomalies* from radiography images (abbreviated as *SQUID*). Radiography imaging protocols focus on particular body regions, therefore producing images of great similarity and yielding recurrent anatomical structures across patients. To exploit this structured information, our *SQUID* consists of a new *Memory Queue* and a novel *in-painting* block in the feature space. We show that *SQUID* can taxonomize the ingrained anatomical structures into recurrent patterns; and in the inference, *SQUID* can identify anomalies (unseen/modified patterns) in the image. *SQUID* surpasses the state of the art in unsupervised anomaly detection by over 5 points on two chest X-ray benchmark datasets. Additionally, we have created a new dataset (*DigitAnatomy*), which synthesizes the spatial correlation and consistent shape in chest anatomy. We hope *DigitAnatomy* can prompt the development, evaluation, and interpretability of anomaly detection methods, particularly for radiography imaging. Code and dataset will be made available.

## 1. Introduction

Vision tasks in photographic imaging and radiography imaging are different. As an example, when identifying objects in photographic images, their locations are generally not important—a cat is a cat no matter if it appears in the left or right of the image. In radiography imaging, on the other hand, the relative location and orientation of a structure are important characteristics that allow the identification of normal anatomy and pathological conditions [16, 82]. Since radiography imaging protocols assess patients in a fairly consistent orientation, the generated images have great similarity across various patients, equipment manufacturers, and facility locations (see examples in Figure 1d). The consistent and recurrent anatomy facilitates the analysis of numerous critical problems and should be considered as a significant advantage for radiography imaging. Several investigations have demonstrated the value of this prior knowledge in

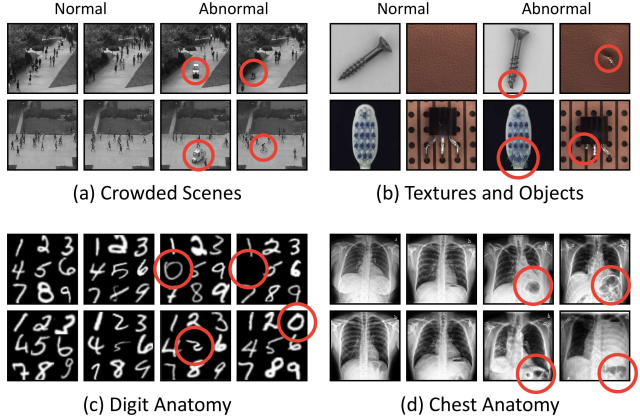


Figure 1. Anomaly detection in radiography images can be both easier and harder than photographic images. It is easier because radiography images are spatially structured due to consistent imaging protocols. It is harder because anomalies in radiography images are subtle and require medical expertise to annotate. With this, we make two contributions: (i) a new *DigitAnatomy* dataset that combines the spatial structure of radiography images and high interpretability of photographic images; (ii) a novel anomaly detection method (*SQUID*) that directly exploits the structured information in radiography images, yielding state-of-the-art performance on two chest radiography benchmarks.

enhancing Deep Nets’ performance by adding location features, modifying objective functions, and constraining coordinates relative to landmarks in images [3, 43, 45, 65, 85]. Our work seeks to answer this critical question: *Can we exploit consistent anatomical patterns and their spatial information to strengthen Deep Nets in detecting anomalies from radiography images without manual annotation?*

We formulate the task of anomaly detection as an in-painting task to exploit this anatomical consistency in appearance, position, and layout. Specifically, we propose *space-aware memory queues* for *in-painting* and *detecting anomalies* from radiography images (abbreviated as *SQUID*). In the training phase, our model can *dynamically* maintain a visual pattern dictionary by taxonomizing recurrent anatomical patterns based on their spatial locations. Due to the consistency in anatomy, the same body

\*Corresponding author: Zongwei Zhou (zzhou82@jh.edu)

region across normal images is expected to express similar visual patterns, which makes the total number of unique patterns manageable. In the inference, since anomaly patterns are not present in the learned dictionary, the generated radiography image is expected to be unrealistic. As a result, the model can identify the anomaly by discriminating the quality of the in-painting task. The success of anomaly detection has two basic assumptions [88]: *first*, anomalies only occur very rarely in the data; *second*, anomalies differ from the normal patterns significantly. Consequently, the resulting dictionary will reflect the general distribution of anatomical patterns in the normal human anatomy.

We have conducted extensive experiments on *two* large-scale, publicly available radiography imaging datasets. Our SQUID is significantly superior to predominant methods in unsupervised anomaly detection by over 5 points on the ZhangLab dataset [29]; remarkably, we have demonstrated a 10-point improvement over the state of the art on the Stanford CheXpert dataset [26]. In addition, we have created a new dataset (DigitAnatomy) to synthesize *spatial correlation* and *consistent shape* of the chest anatomy in radiographs. DigitAnatomy was created to prompt the development, evaluation, and interpretability of anomaly detection methods. The qualitative visualization clearly shows the superiority of our SQUID over the current state of the art.

In summary, our **contributions** include: (i) the best performing unsupervised anomaly detection method for chest radiography imaging; (ii) a new dataset to promote anomaly detection research. Technically, our SQUID overcomes limitations in dominant unsupervised anomaly detection methods [1, 13, 32, 57, 81] by inventing Memory Queue (§3.2), and Feature-level In-painting (§3.3).

## 2. Related Work

### 2.1. Anomaly Detection in Natural Imaging

Anomaly detection is the task of identifying rare events that deviate from the distribution of normal data [48]. Early attempts include one-class SVM [60], dictionary learning [80], and sparse coding [7]. Due to the lack of sufficient samples of anomalies, later works typically formulate anomaly detection as an unsupervised learning problem [8, 22, 23, 34, 35, 39, 53, 63, 89]. These can be roughly categorized into reconstruction-based and density-based methods. Reconstruction-based methods train a model (e.g. Auto-Encoder) to recover the original inputs [6, 62, 67, 86, 87]. The anomalies are identified by subtracting the reconstructed image from the input image. Density-based methods predict anomalies by estimating the normal data distribution (e.g. via VAEs [32] or GANs [2, 57, 58]). However, the learned distribution for normal images by these methods cannot explain the possible abnormalities. In this paper, we address these limitations by maintaining a visual pattern dictionary

which is extracted from homogeneous medical images.

Several other previous works investigated the use of image in-painting for anomaly detection, *i.e.* parts of the input image are masked out and the model is trained to recover the missing parts in a self-supervised way [19, 37, 47, 52, 78]. There are also plenty of works on detecting anomalies in video sequences [10, 18, 41, 42]. Recently, Bergmann *et al.* [4] and Salehi *et al.* [55] proposed student-teacher networks similar to ours, whereas our method utilizes such a structure to distillate input-aware features only, and the teacher network is completely disabled during inference.

### 2.2. Anomaly Detection in Medical Imaging

Compared with natural imaging, anomaly detection in medical imaging is relatively sparse [12, 69] and most of them heavily rely on manual annotation [59]. With the help of GANs, anomaly detection can be achieved with *weak* annotation. In AnoGAN [58], the discriminator was heavily over-fitted to the normal image distribution to detect the anomaly. Subsequently, f-AnoGAN [57] was proposed to improve the computational efficiency. Marimont *et al.* [46] designed an auto-decoder network to fit the distribution of normal images. The spatial coordinates and anomaly probabilities are mapped over a proxy for different tissue types. Han *et al.* [17] proposed a two-step GAN-based framework for detecting anomalies in MRI slices as well. However, their method relies on a voxel-wise representation for the 3D MRI sequences, which is impossible in our task. Most recently, a hybrid framework SALAD [81] was proposed that combines GAN with self-supervised techniques. Normal images are first augmented to carry the forged anomaly through pixel corruption and pixel shuffling. The fake abnormal images, along with the original normal ones, are fed to the GAN for learning more robust feature representations. However, these approaches demand strong prior knowledge and assumptions about the anomaly type to make the augmentation effective.

Differing from photographic images, radiography imaging protocols produce images with consistent anatomical patterns, and meanwhile, the anomalies in radiography images can be subtle in appearance and hard to interpret (Figure 1). Different from most existing works, we present a novel method that explicitly harnesses the radiography images’ properties, therefore dramatically improving the performance in anomaly detection from radiography images.

### 2.3. Memory Networks

Incorporating memory modules into neural networks has been demonstrated to be effective for many tasks [5, 11, 28, 33, 36]. Adopting Memory Matrix for unsupervised anomaly detection was first proposed in MemAE [13]. In addition to auto-encoding (AE), Gong *et al.* injected an extra Memory Matrix between the encoder and the decoder

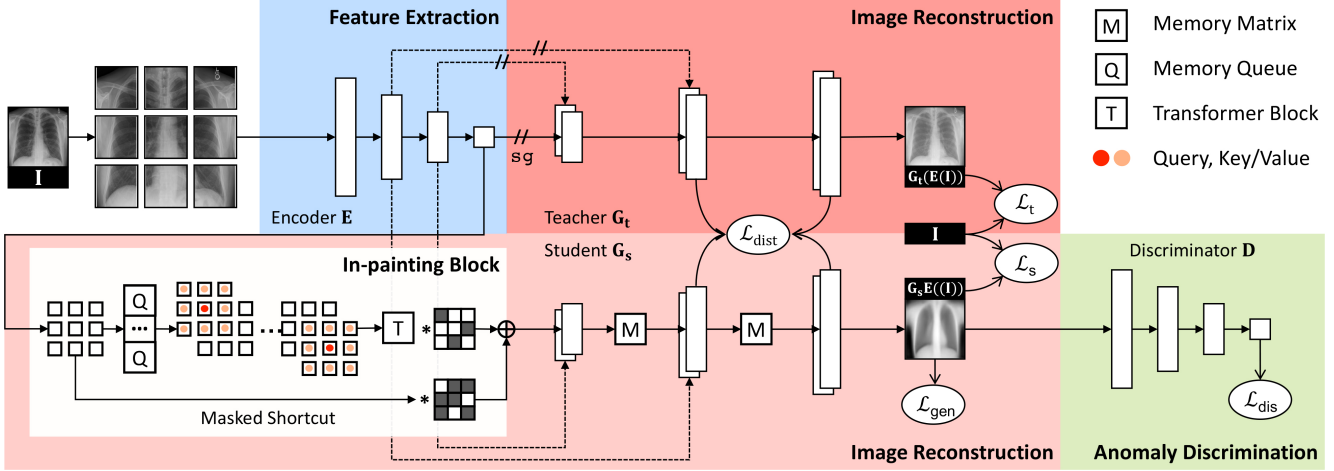


Figure 2. **SQUID**. We divide an input image into  $N \times N$  non-overlapping patches and feed them into the encoder for feature extraction. Two generators will be trained to reconstruct the original image. Along with the reconstruction, a dictionary of anatomical patterns will be created and updated dynamically via a novel Memory Queue (§3.2); The teacher generator directly uses the features extracted by the encoder; the student generator uses the features augmented by our new in-painting block (§3.3). The teacher and student generators are coupled through a knowledge distillation paradigm. We employ a discriminator to assess whether the image reconstructed by the student generator is real or fake. Once trained, the discriminator can be used to detect anomalies in test images (§3.4).

to capture normal feature patterns during training. The matrix is jointly optimized along with the AE and hence learns an essential basis to be able to assemble normal patterns. Based on this paradigm, Park *et al.* [49] introduced a non-learnable memory module that can be updated with inputs. Note that although our proposed Memory Queue also does not require any gradients, our method differs significantly in its usage purpose and updating rules. Considering the extra memory usage in existing methods, Lv *et al.* [44] proposed a dynamic prototype unit that encodes normal dynamics on the fly, while consuming little additional memory. In this paper, we overcome the limitations of Memory Matrix and propose the effective yet efficient Memory Queue for unsupervised anomaly detection in radiography images.

### 3. SQUID

#### 3.1. Overview

**Feature Extraction:** We divide the input image into  $N \times N$  non-overlapping patches and feed them into an encoder for feature extraction. The extracted features will be used for image reconstruction. Practically, the encoder can be any backbone architectures [9, 66], and for simplicity, we adopt basic Convolutions and Pooling layers in this work.

**Image Reconstruction:** We introduce teacher and student generators to reconstruct the original image. Along with the reconstruction, a dictionary of anatomical patterns will be created and updated dynamically via a novel Memory Queue (§3.2). Specifically, the teacher generator reconstructs the image using the features extracted by the encoder

directly (essentially an auto-encoder [54]). The student generator, on the other hand, reconstructs the image using the features augmented by our in-painting block (§3.3). The teacher and student generators are coupled through a knowledge distillation paradigm [24] at all the up-sampling levels. The objective of the student generator is to reconstruct a normal image from the augmented features; the reconstructed image will then be used for anomaly discrimination (§3.4); while the teacher generator<sup>1</sup> serves as a regularizer to prevent the student from collapsing (constantly generating the same normal image).

**Anomaly Discrimination:** Following the adversarial learning [57, 58], we employ a discriminator to assess whether the generated image is real or fake. Both teacher and student generators will receive the gradient derived from the discriminator. The two generators and the discriminator are competing against each other in a way that, together, they converge to an equilibrium. Once trained, the discriminator can be used to detect anomalies in test images (§3.4).

#### 3.2. Inventing Memory Queue as Dictionary

**Motivation:** The Memory Matrix was initially introduced by Gong *et al.* [13] and has since been widely adopted in unsupervised anomaly detection [14, 30, 41, 77, 83]. To forge a “normal” appearance, the features are *augmented* by weighted averaging the similar patterns in Memory Matrix. This augmentation is, however, applied to the features

<sup>1</sup>We disabled the backpropagation between the teacher and encoder by stop-gradient [20] and showed its empirical benefit in Table 2.

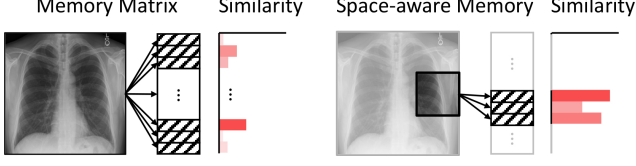


Figure 3. **Space-aware memory.** For unique encoding of location information, we restrict each patch to be only accessible by a non-overlapping region in the memory.

extracted from the whole image, discarding the location and spatial information embedded in images. Therefore, Memory Matrix in its current form cannot perceive the anatomical consistency that radiography images can offer.

**Space-aware Memory:** To harness the spatial information, we pass the divided small patches into the model rather than the whole image. These patches are associated with unique location information of the original image. We seek to build the relationship between the patch location and memory region. In doing so, we narrow the search space in Memory Matrix subject to the patch location of the original image. That is, a patch derived from a particular location can only search for similar patterns within a specific region in Memory Matrix (illustrated in Figure 3). We refer to this new strategy as “space-aware memory” because it enables encoding of the spatial information into Memory Matrix. Space-aware memory can also accelerate the searching speed compared with [13] as it no longer has to go through the entire Memory Matrix to assemble similar features.

**Memory Queue:** Figure 4 indicates that the learned patterns in Memory Matrix [13] (blue dots) distribute differently from the patch features of the training set (gray dots). It is caused by the learning-based nature of Memory Matrix, which forges *fake normal* patterns by combining the learned basis. The discrepancy between reassembled normal features and the patch features makes it hard for the generator to reconstruct the normal image. To address this issue, we propose a Memory Queue to interpolate amongst *real normal* patterns that exist in the training set, therefore presenting a consistent distribution of the patch features (red dots). Specifically, it directly copies previously seen features into Memory Queue during training<sup>2</sup>. Once trained, Memory Queue can be used as a *dictionary* of normal anatomical patterns. The patterns in the dictionary will be assembled weighted by their similarity to the patch feature.

<sup>2</sup>In practice, copying features into Memory Queue at every training iteration demands considerable computational time. Supposing  $N$  patterns in Memory Queue and  $M$  training iterations, the sampling strategy in [13] demands an  $\mathcal{O}(NM)$  time complexity. We implement it more efficiently: at each iteration, the current batch of features will be copied into Memory Queue for *only once* (Figure 5c), yielding a linear complexity of  $\mathcal{O}(M + cM)$  with the copy-and-paste operation in a constant time  $c$ . We follow the first-in-last-out (FILO) paradigm to update Memory Queue continuously.

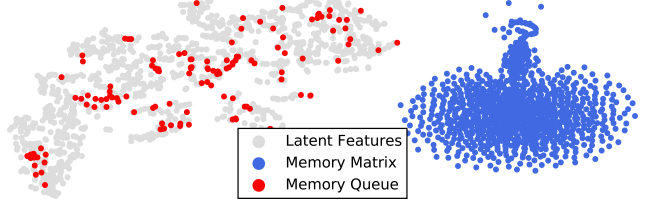


Figure 4. t-SNE visualizations of patterns in Memory Matrix, Memory Queue, and patch features of the training samples [70]. Patterns in Memory Matrix are far away from the distribution of patch features, while patterns in Memory Queue (as copies of previously seen features) share a similar distribution.

**Gumbel Shrinkage:** Controlling the number of activated patterns in the memory has proven to be advantageous for anomaly detection [13, 15]. However, setting a hard shrinkage threshold fails to adapt to cases where abnormal signals are sufficient to reconstruct a normal image. Therefore, we consider a soft threshold to activate the top- $k$  most similar patterns in Memory Queue, wherein only the top- $k$  patterns will receive gradients, and the remainder could not be updated as expected. To extend gradient descent to more patterns in Memory Queue, inspired by Jang [27], we present a *Gumbel Shrinkage* schema: only activating the top- $k$  most similar patterns during the forward pass and distributing the gradient to all patterns during back-propagation. Our Gumbel Shrinkage schema is formulated as:

$$\mathbf{w}' = \text{sg}(\text{hs}(\mathbf{w}, \text{topk}(\mathbf{w})) - \phi(\mathbf{w})) + \phi(\mathbf{w}), \quad (1)$$

where  $\mathbf{w}$  denotes the similarity between the patch feature and patterns in Memory,  $\text{sg}(\cdot)$  the stop-gradient operation,  $\text{hs}(\cdot, t)$  the hard shrinkage operator with threshold  $t$ , and  $\phi(\cdot)$  the Softmax function. In the forward pass, Gumbel Shrinkage ensures the combination of top- $k$  most similar patterns in Memory Queue; During the back-propagation, Gumbel Shrinkage is essentially a Softmax function. We apply Gumbel Shrinkage to Memory Queue in the in-painting block and Memory Matrix in the student generator.

### 3.3. Formulating Anomaly Detection as In-painting

**Motivation:** Image in-painting [38, 50] was initially proposed to recover corrupted regions in the image based on the available neighboring context. The recovered regions, however, have been seen to associate with boundary artifacts, distorted and blurry predictions, particularly when using methods based on Deep Nets [40]. These undesired artifacts are responsible for numerous false positives when formulating anomaly detection as an image in-painting task [62, 86]. It is because the subtraction between input and output will reveal artifacts generated by Deep Nets instead of true anomalies. To alleviate this issue, we propose the in-painting task at the feature level rather



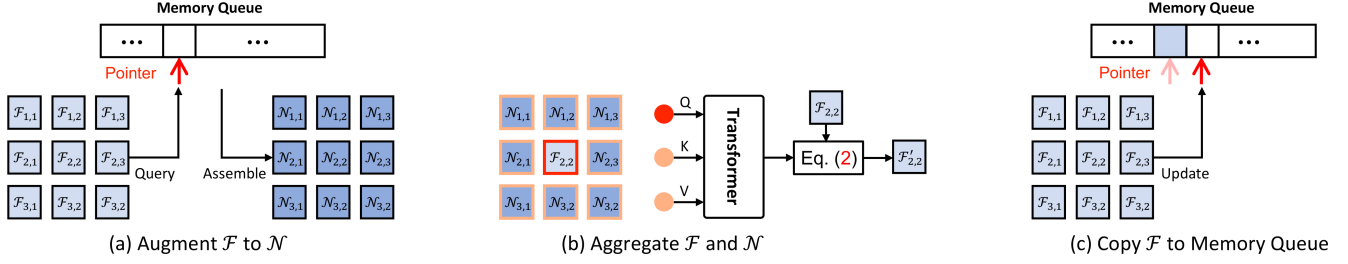


Figure 5. **Three-step workflow of our in-painting block.** (a) Each non-overlapping patch feature  $\mathcal{F}$  is queried to an unique region in Memory Queue, the most similar items are assembled to  $\mathcal{N}$ . (b) Each center patch feature  $\mathcal{F}$  and its eight neighbors  $\mathcal{N}$  are used as query and key/value respectively to a Transformer layer for in-painting. (c) Each specific Memory Queue region copies their corresponding patch features  $\mathcal{F}$  into the memory by maintaining a pointer. Note that this step is only performed during training.

than the image pixel level. Latent features are invariant to subtle noise, rotation, and translation in the pixel level and therefore are expected to be more suitable for anomaly detection. Following conventional in-painting intuition, the model predicts central features based on neighboring features and augments patch features in a sliding-window style.

**In-painting Block:** We integrate our Memory Queue with a novel in-painting block to perform an in-painting task at the feature level. The  $w \times h$  non-overlapping patch features  $\mathcal{F}_{\{(1,1), \dots, (w,h)\}}$  are augmented to the most similar “normal” patterns  $\mathcal{N}_{\{(1,1), \dots, (w,h)\}}$  in Memory Queue (Figure 5a). Since  $\mathcal{N}$  is assembled by patterns from previously seen images, it is not subject to the current input image. To recap characteristics of the current image, naturally, we aggregate both patch features  $\mathcal{F}$  and augmented features  $\mathcal{N}$  using a Transformer block (one multi-head attention layer + one MLP) [71]. For each patch  $\mathcal{F}_{i,j}$ , its spatially adjacent eight “normal” patches  $\mathcal{N}_{\{(i-1,j-1), \dots, (i+1,j+1)\}}$  are used as conditions to refine  $\mathcal{F}_{i,j}$  (Figure 5b). The query token is flattened  $\mathcal{F}_{(i,j)} \in \mathcal{R}^{1 \times *}$  and key/value tokens<sup>3</sup> are  $\mathcal{N}_{\{(i-1,j-1), \dots, (i+1,j+1)\}} \in \mathcal{R}^{8 \times *}$ . At the start and the end of our in-painting block, we apply an extra pair of point-wise convolutions ( $1 \times 1$  convolutional kernel) [21] to reduce feature dimensions and accelerate the training process.

**Masked Shortcut:** We employ a shortcut within the in-painting block to aggregate features and ease optimization. Our empirical study shows that a direct residual connection can downgrade the effectiveness of the in-painting block (detailed in Appendix C). Inspired by Xiang *et al.* [76], we utilize a random binary mask to gate feature shortcuts during training (Figure 5b). As such, given the features  $\mathcal{F}$ , the output of the in-painting block is obtained by:

$$\mathcal{F}' = (1 - \delta) \cdot \mathcal{F} + \delta \cdot \text{inpaint}(\mathcal{F}), \quad (2)$$

where  $\text{inpaint}(\cdot)$  is the in-painting block described earlier,  $\delta \sim \text{Bernoulli}(\rho)$  is a binary variable and  $\rho$  is a hyper-

<sup>3</sup>We use zero padding for out of range patches.

parameter to control the gating probability. During inference, we disable the shortcut completely such that  $\mathcal{F}' = \text{inpaint}(\mathcal{F})$  for deterministic predictions.

### 3.4. Anomaly Discrimination

Our in-painting block focuses on augmenting any patch feature (either normal or abnormal) into the normal feature pattern. The student generator will then reconstruct a “normal” image based on the augmented features. The teacher generator is used to preserve the normal image intact and prevent the student from collapsing. Once trained, the semantic (rather than pixel-level) difference between the input and the reconstructed image is expected to be small if normal; the semantic difference will be big if there are anomalies. We therefore delegate the optimized discriminator network for alerting anomalies perceptually. For better clarification, we notate the encoder, teacher generator, student generator, and discriminator as  $\mathbf{E}$ ,  $\mathbf{G}_t$ ,  $\mathbf{G}_s$ , and  $\mathbf{D}$ . An anomaly score ( $A$ ) can be computed through:  $A = \phi\left(\frac{\mathbf{D}(\mathbf{G}_s(\mathbf{E}(\mathbf{I}))) - \mu}{\sigma}\right)$ , where  $\phi(\cdot)$  is the Sigmoid function,  $\mu$  and  $\sigma$  are the mean and standard deviation of anomaly scores calculated on all training samples.

### 3.5. Loss Function

Our SQUID is optimized by five loss functions. The mean square error (MSE) between input and reconstructed images is used for both teacher and student generators. Concretely,  $\mathcal{L}_t = (\mathbf{I} - \mathbf{G}_t(\mathbf{E}(\mathbf{I})))^2$  and  $\mathcal{L}_s = (\mathbf{I} - \mathbf{G}_s(\mathbf{E}(\mathbf{I})))^2$  for the teacher and student generators, respectively, where  $\mathbf{I}$  denotes the input image. Following the knowledge distillation paradigm, we apply a distance constraint between the teacher and student generators to all levels of features:  $\mathcal{L}_{\text{dist}} = \sum_{i=1}^l (\mathcal{F}_t^i - \mathcal{F}_s^i)^2$ , where  $l$  is the level of features used for knowledge distillation,  $\mathcal{F}_t$  and  $\mathcal{F}_s$  are the intermediate features in the teacher and student generators, respectively. In addition, we employ an adversarial loss (similar to DCGAN [51]) to improve the quality of the image generated by the student gener-

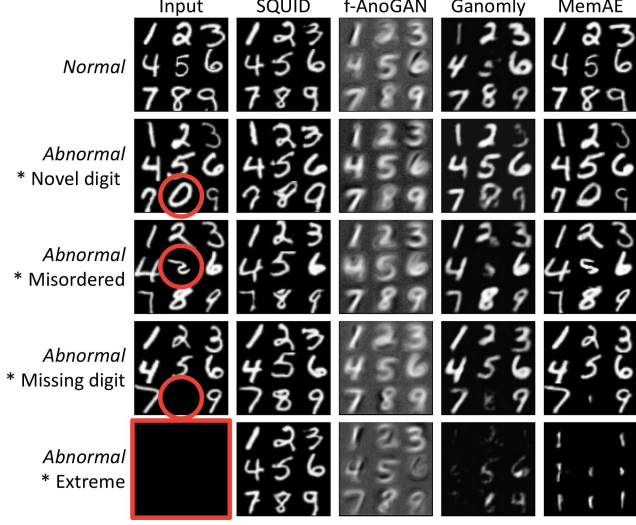


Figure 6. Comparisons of reconstruction results on DigitAnatomy of our SQUID, f-AnoGAN [57], Ganomaly [1], and MemAE [13]. More visualization can be found in Appendix F.

ator. Specifically, the following equation is minimized:  $\mathcal{L}_{\text{gen}} = \log(1 - D(G_s(E(I))))$ . The discriminator seeks to maximize the average of the probability for real images and the inverted probability for fake images:  $\mathcal{L}_{\text{dis}} = \log(D(I)) + \log(1 - D(G_s(E(I))))$ .

In summary, our SQUID is trained to *minimize* the generative loss terms ( $\lambda_t \mathcal{L}_t + \lambda_s \mathcal{L}_s + \lambda_{\text{dist}} \mathcal{L}_{\text{dist}} + \lambda_{\text{gen}} \mathcal{L}_{\text{gen}}$ ) and to *maximize* the discriminative loss term ( $\lambda_{\text{dis}} \mathcal{L}_{\text{dis}}$ ).

## 4. Experiments

### 4.1. New Benchmark

**DigitAnatomy:** We have created a synthetic dataset to verify the main idea, wherein the human anatomy is translated into Arabic digits one to nine in an in-grid placement (see examples in Figure 1 and Figure 6). The images containing digits one to nine in the correct order are considered “normal”; otherwise, as “abnormal”. The types of simulated abnormalities include missing, misordered, and flipped digit(s) and the digit zero. We provide the pseudocode for creating DigitAnatomy in Appendix E. This DigitAnatomy dataset is particularly advantageous for radiography imaging for the following three reasons. *First*, it simulates two unique properties of radiography images, *i.e.* spatial correlation and consistent shape. *Second*, annotating radiography images demands specialized expertise, but digits are easier for problem shooting and debugging. *Third*, the ground truth of the simulated anomaly is readily accessible in DigitAnatomy, whereas it is hard to collect sufficient examples for each abnormal type in radiography images.

### 4.2. Public Benchmarks

**ZhangLab Chest X-ray [29]:** This dataset contains healthy and pneumonia (as anomaly) images, *officially* split into training and test sets. The training set consists of 1,349 normal and 3,883 abnormal images; the test set has 234 normal and 390 abnormal images. We randomly separate 200 images (100 normal and 100 abnormal) from the training set as the validation set for early-stopping. Since the images are of varying sizes, we resized all the images to  $128 \times 128$  without losing essential details. We used this dataset for ablation studies as well.

**Stanford CheXpert [26]:** We conducted evaluations on the front-view PA images in the CheXpert dataset, which account for a total of 12 different anomalies. In all front-view PA images, there are 5,249 normal and 23,671 abnormal images for training; 250 normal and 250 abnormal images (with at least 10 images per disease type) from the training set for testing; 14 normal and 19 abnormal images for early-stopping (val set based on the *official* split). All images are resized to  $128 \times 128$  as inputs.

### 4.3. Baselines and Metrics

There are not many unsupervised anomaly detection baselines available for direct comparison due to three reasons: (i) no publicly released code [25, 72, 74, 81], (ii) requiring weak/strong manual annotation [56, 68, 75, 84], and (iii) relying on additional information from other modalities [17, 46]. Therefore, we considered six major baselines for direct comparison: MemAE [13]—a Memory Matrix based method, Ganomaly [1]—a GAN-based method, f-AnoGAN [57]—the current state of the art for medical imaging, and CutPaste [37], PANDA [52], M-KD [56]—the most recent unsupervised anomaly detection methods. We also compared our SQUID with SALAD [81] using the results reported in their paper in the same test set. We evaluated performance using standard metrics: Receiver Operating Characteristic (ROC) curve, Area Under the ROC Curve (AUC), Accuracy (Acc), and F1-score (F1). All results were based on at least *three* independent runs.

### 4.4. Implementation Details

We utilized common data augmentation strategies such as random translation within the range  $[-0.05, +0.05]$  in four directions and a random scaling within the range of  $[0.95, 1.05]$ . The Adam [31] optimizer was used with a batch size of 16 and a weight decay of  $1e-5$ . The learning rate was initially set to  $1e-4$  for both the generator and the discriminator and then decayed to  $2e-5$  in 200 epochs following the cosine annealing scheduler. The discriminator is trained at every iteration, while the generator is trained every two iterations. We set loss weights as  $\lambda_t = 0.01$ ,  $\lambda_s = 10$ ,  $\lambda_{\text{dist}} = 0.001$ ,  $\lambda_{\text{gen}} = 0.005$ , and  $\lambda_{\text{dis}} = 0.005$ . We

<i>ZhangLab</i>	Ref & Year	AUC (%)	Acc (%)	F1 (%)
Auto-encoder <sup>†</sup>	-	59.9	63.4	77.2
VAE <sup>†</sup> [32]	Arxiv'13	61.8	64.0	77.4
Ganomaly <sup>†</sup> [1]	ACCV'18	78.0	70.0	79.0
f-AnoGAN <sup>†</sup> [57]	MIA'19	75.5	74.0	81.0
MemAE [13]	ICCV'19	77.8±1.4	56.5±1.1	82.6±0.9
SALAD <sup>†</sup> [81]	TMI'21	82.7±0.8	75.9±0.9	82.1±0.3
CutPaste [37]	CVPR'21	73.6±3.9	64.0±6.5	72.3±8.9
PANDA [52]	CVPR'21	65.7±1.3	65.4±1.9	66.3±1.2
M-KD [55]	CVPR'21	74.1±2.6	69.1±0.2	62.3±8.4
SQUID	-	<b>87.6±1.5</b>	<b>80.3±1.3</b>	<b>84.7±0.8</b>

<sup>†</sup>The results are taken from Zhao *et al.* [81]

<i>CheXpert</i>	Ref & Year	AUC (%)	Acc (%)	F1 (%)
Ganomaly [1]	ACCV'18	68.9±1.4	65.7±0.2	65.1±1.9
f-AnoGAN [57]	MIA'19	65.8±3.3	63.7±1.8	59.4±3.8
MemAE [13]	ICCV'19	54.3±4.0	55.6±1.4	53.3±7.0
CutPaste [37]	CVPR'21	65.5±2.2	62.7±2.0	60.3±4.6
PANDA [52]	CVPR'21	68.6±0.9	66.4±2.8	65.3±1.5
M-KD [55]	CVPR'21	69.8±1.6	66.0±2.5	63.6±5.7
SQUID	-	<b>78.1±5.1</b>	<b>71.9±3.8</b>	<b>75.9±5.7</b>

Table 1. Benchmark results on the test sets of two datasets.

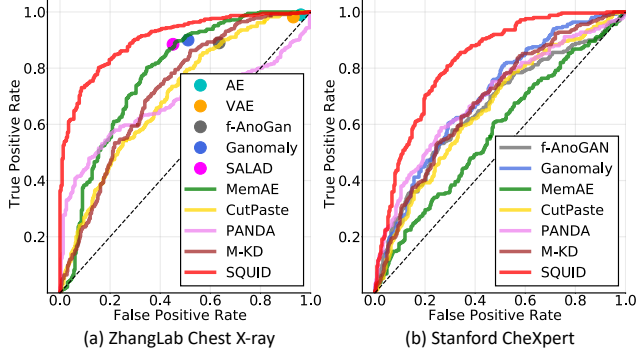


Figure 7. ROC curves comparison on the two datasets.

divide the input images in  $2 \times 2$  non-overlapping patches, fix the shortcut mask probability at  $p = 95\%$ , and activate only the top 5 similar patterns in the Gumbel Shrinkage. The choice of these hyper-parameters is studied in §5.3. The architectures of our generators and discriminator are detailed in Appendix B. The pseudocode of SQUID is provided in Appendix A. Code and dataset will be available.

## 5. Results

### 5.1. Interpreting SQUID on DigitAnatomy

Figure 6 presents qualitative results on DigitAnatomy to examine the capability of image reconstruction and to interpret the mistakes made by existing methods [1, 13, 57]. We deliberately inject anomalies (*e.g.* novel, misordered, missing digits) into normal images (highlighted in red) and test if the model can reconstruct their normal counterparts. We also assess the reconstruction quality from a blank image (as an extreme case) to raise the task difficulty. In gen-

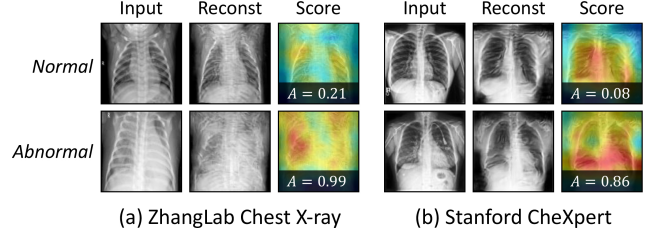


Figure 8. Reconstruction results of SQUID on the two datasets. The corresponding anomaly scores are shown at the bottom. More visualization can be found in Appendix F.

eral, the images reconstructed by our SQUID carry more meaningful and indicative information than other baseline methods. It is mainly attributed to our *space-aware* memory, with which the resulting dictionary is associated with unique patterns as well as their spatial information. Once an anomaly arises (*e.g.* missing digit), the in-painting block will augment the abnormal feature to its normal counterpart by assembling top- $k$  most similar patterns from the dictionary. Other methods, however, do not possess this ability, so they reconstruct defective images. For instance, GAN-based methods (f-AnoGAN and Ganomaly) tend to reconstruct an exemplar image averaged from the training examples. MemAE performs relatively better due to its Memory Matrix, but it does not work well for the anomaly of missing digits and completely fails on the extreme anomaly attack.

### 5.2. Benchmarking SQUID on Chest X-rays

Our SQUID was mainly evaluated on two large-scale benchmarks: ZhangLab Chest X-ray and Stanford CheXpert for comparing with a wide range of state-of-the-art counterparts. According to Table 1, our SQUID achieves the most promising result in terms of all metrics on both datasets. Specifically, SQUID outperforms the second best runner-up counterparts by at least 5% in AUC, 5% in Accuracy. The highest F1 scores SQUID achieved, along with the ROC curves shown in Figure 7, demonstrate that our method yields the best trade-off between sensitivity and specificity. Overall, the significant improvements observed with SQUID proved the effectiveness of our proposed designs and techniques in this work.

In Figure 8, we visualize the reconstruction results of SQUID on exemplary normal and abnormal images in the two datasets. For normal cases, SQUID can easily find a similar match in Memory Queue and hence achieves the reconstruction smoothly. For abnormal cases, contradiction will arise by imposing forged normal patterns into the abnormal features. In this way, the generated images will vary significantly from the input, which will then be captured by the discriminator. We plot the heatmap of the discriminator (using Grad-CAM [61]) to indicate the most likely regions to appear anomalous. As a result, the reconstructed healthy

Method	AUC (%)	Acc (%)	F1 (%)
w/o Space-aware Memory	77.6±0.5	75.5±0.5	82.5±0.6
w/o In-painting Block	80.9±2.1	75.8±1.5	81.6±1.3
w/o Gumbel Shrinkage	81.1±0.9	77.6±0.9	81.3±0.8
w/o Knowledge Distillation	81.2±0.8	75.2±0.7	81.3±0.8
w/o Stop Gradient	81.7±4.3	76.7±2.8	82.5±1.6
w/o Memory Queue	82.5±1.1	78.6±0.9	81.7±1.1
w/o Masked Shortcuts	82.5±1.3	76.4±0.8	82.3±1.1
w/o Decoder Memory	82.9±1.2	77.4±1.1	81.2±0.5
Full SQUID	<b>87.6±1.5</b>	<b>80.3±1.3</b>	<b>84.7±0.8</b>

Table 2. Component studies indicate that the overall performance benefits from all the components in SQUID.

images yield much lower anomaly scores than the diseased ones, validating the effectiveness of SQUID.

**Limitation:** We found SQUID in its current form, is not able to *localize* anomalies at the pixel level precisely. It is understandable because unlike [56, 64, 68, 73, 79], our SQUID is an unsupervised method, requiring zero manual annotation for normal/abnormal images. More investigation on pixel-level localization (or even segmentation) and multi-scale detections could be meaningful in the future.

### 5.3. Ablating Key Properties in SQUID

**Component Study:** We examine the impact of components in SQUID by taking each one of them out of the entire framework. Table 2 shows that each component accounts for at least 5% performance gain. The space-aware memory (+10.0%) and in-painting block (+6.7%) are the top-2 most significant contributors, which underline our motivation and justification of the method development (§3.2 and §3.3). Moreover, the knowledge distillation from teacher to student generators strikes an important balance: the student generator reconstructs faithful “normal” images from similar anatomical patterns in the dictionary while preserving unique characteristics of each input image (regularized by the teacher generator). Besides, we must acknowledge that the training tricks (*e.g.* hard shrinkage [27], stop gradient [20], masked shortcut [76]) are necessary for the remarkable performance. Although replacing Memory Queue with Memory Matrix could maintain a decent result (only dropped 5.1%), our Memory Queue presents a more trustworthy recovery of “normal” patterns in the image than Memory Matrix (MemAE [13]), evidenced by Figure 6.

**Hyper-parameter Robustness:** The number of patch divisions, the topk value in Gumbel Shrinkage, the number of memory patterns within a specific region of Memory Queue, and the shortcut masking probability  $\rho$  are four important hyper-parameters of SQUID. Here, we conducted exhausted experiments on these parameters in Figure 9. Trials were first made on the number of patches from  $1 \times 1$  to  $8 \times 8$ . When dividing input images into a single patch, space-aware settings are not triggered, hence yielding the worst performance. Although the spatial structures are rel-

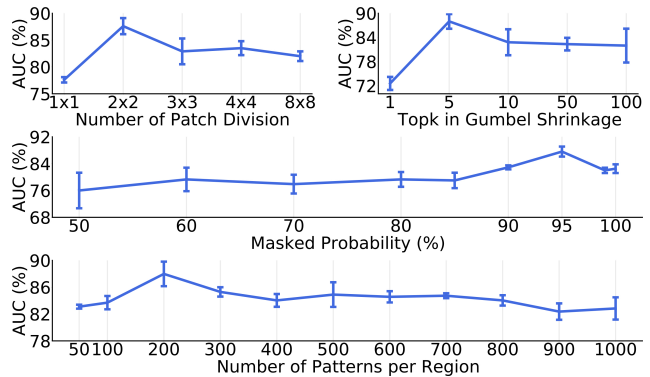


Figure 9. SQUID is robust to hyper-parameter modifications. The best result is obtained at dividing  $2 \times 2$  patches, setting 200 patterns per memory region, using 95% mask probability, and activating top 5 patterns through Gumbel Shrinkage. Tabular results are in Appendix D.

atively stable in most chest X-rays, certain deviations can still be observed. Therefore, with small patches, object parts in one patch can easily appear in adjacent patches and be misdetected as anomalies. The number of topk activations in Gumbel softmax also impacts the performances. By assembling the top-5 most similar patterns through Gumbel softmax, SQUID is able to achieve the best result. When replacing input features with the top-1 most similar pattern, SQUID suffers from a performance drop by -15% AUC. According to the AUC vs. number of patterns in each Memory Queue region, we found that a small number of items is sufficient to support normal pattern querying in local regions and the best result is achieved by using merely 200 items per region. When the item number exceeds 500 per region, AUC scores begin to drop continuously. AUC vs. the mask probability  $\rho$  was further plotted to verify that enabling a limited number of feature skips ( $\rho = 95\%$ ) yields the best AUC score. The effectiveness of the in-painting will severely deteriorate if more features are allowed to be skipped ( $\rho < 90\%$ ).

## 6. Conclusion

We present SQUID for unsupervised anomaly detection from radiography images. Qualitatively, we show that SQUID can taxonomize the ingrained anatomical structures into recurrent patterns; and in the inference, SQUID can identify anomalies (unseen/modified patterns) in the image. Quantitatively, SQUID is significantly superior to predominant methods in unsupervised anomaly detection by over 5 points on the ZhangLab dataset; remarkably, SQUID achieves a 10-point improvement over the state of the art on the Stanford CheXpert dataset. The competitive results are attributable to our observation: *Radiography imaging protocols focus on particular body regions, therefore producing*



images of great similarity and yielding recurrent anatomical structures across patients. Additionally, we created a new dataset (DigitAnatomy) that synthesizes the spatial correlation and consistent shape of chest anatomy in radiography images. DigitAnatomy was created to prompt the development, evaluation, and interpretability of anomaly detection methods, particularly for radiography imaging.

## Acknowledgments

This work was supported by the Lustgarten Foundation for Pancreatic Cancer Research. We appreciate the constructive suggestions from Yingda Xia, Yixiao Zhang, Bowen Li, Adam Kortylewski, and Huiyu Wang.

## References

- [1] Samet Akcay, Amir Atapour-Abarghouei, and Toby P Breckon. Ganomaly: Semi-supervised anomaly detection via adversarial training. In *Asian conference on computer vision*, pages 622–637. Springer, 2018. 2, 6, 7, 16, 18
- [2] Varghese Alex, Mohammed Safwan KP, Sai Saketh Chen-namsetty, and Ganapathy Krishnamurthi. Generative adversarial networks for brain lesion detection. In *Image Processing Medical Imaging*, volume 10133, page 101330G. International Society for Optics and Photonics, 2017. 2
- [3] Emran Mohammad Abu Anas, Abtin Rasoulia, Alexander Seitel, Kathryn Darras, David Wilson, Paul St John, David Pichora, Parvin Mousavi, Robert Rohling, and Purang Abol-maesumi. Automatic segmentation of wrist bones in ct using a statistical wrist shape + pose model. *IEEE transactions on medical imaging*, 35(8):1789–1801, 2016. 1
- [4] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. Uninformed students: Student-teacher anomaly detection with discriminative latent embeddings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4183–4192, 2020. 2
- [5] Qi Cai, Yingwei Pan, Ting Yao, Chenggang Yan, and Tao Mei. Memory matching networks for one-shot image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4080–4088, 2018. 2
- [6] Xiaoran Chen and Ender Konukoglu. Unsupervised detection of lesions in brain mri using constrained adversarial auto-encoders. *Medical Imaging with Deep Learning*, 2018. 2
- [7] Yang Cong, Junsong Yuan, and Ji Liu. Sparse reconstruction cost for abnormal event detection. In *CVPR 2011*, pages 3449–3456. IEEE, 2011. 2
- [8] Terrance DeVries and Graham W Taylor. Learning confidence for out-of-distribution detection in neural networks. *arXiv preprint arXiv:1802.04865*, 2018. 2
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *International Conference on Learning Representations*, 2020. 3
- [10] Patrick Esser, Ekaterina Sutter, and Björn Ommer. A variational u-net for conditional appearance and shape generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8857–8866, 2018. 2
- [11] Chenyou Fan, Xiaofan Zhang, Shu Zhang, Wensheng Wang, Chi Zhang, and Heng Huang. Heterogeneous memory enhanced multimodal attention model for video question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1999–2007, 2019. 2
- [12] Tharindu Fernando, Harshala Gammulle, Simon Denman, Sridha Sridharan, and Clinton Fookes. Deep learning for medical anomaly detection—a survey. *arXiv preprint arXiv:2012.02364*, 2020. 2
- [13] Dong Gong, Lingqiao Liu, Vuong Le, Budhaditya Saha, Moussa Reda Mansour, Svetha Venkatesh, and Anton

- van den Hengel. Memorizing normality to detect anomaly: Memory-augmented deep autoencoder for unsupervised anomaly detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1705–1714, 2019. 2, 3, 4, 6, 7, 8, 16, 18
- [14] Dong Gong, Zhen Zhang, Javen Qinfeng Shi, and Anton van den Hengel. Memory-augmented dynamic neural relational inference. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11843–11852, 2021. 3
- [15] Alex Graves, Greg Wayne, and Ivo Danihelka. Neural Turing machines. *arXiv preprint arXiv:1410.5401*, 2014. 4
- [16] Fatemeh Haghighi, Mohammad Reza Hosseinzadeh Taher, Zongwei Zhou, Michael B Gotway, and Jianming Liang. Transferable visual words: Exploiting the semantics of anatomical patterns for self-supervised learning. *IEEE Transactions on Medical Imaging*, 2021. 1
- [17] Changhee Han, Leonardo Rundo, Kohei Murao, Tomoyuki Noguchi, Yuki Shimahara, Zoltán Ádám Milacski, Saori Koshino, Evis Sala, Hideki Nakayama, and Shin’ichi Satoh. Madgan: unsupervised medical anomaly detection gan using multiple adjacent brain mri slice reconstruction. *BMC bioinformatics*, 22(2):1–20, 2021. 2, 6
- [18] Mahmudul Hasan, Jonghyun Choi, Jan Neumann, Amit K Roy-Chowdhury, and Larry S Davis. Learning temporal regularity in video sequences. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 733–742, 2016. 2
- [19] Matthias Haselmann, Dieter P Gruber, and Paul Tabatabai. Anomaly detection using deep learning based image completion. In *2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA)*, pages 1237–1242. IEEE, 2018. 2
- [20] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020. 3, 8
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 5
- [22] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *International Conference on Learning Representations*, 2016. 2
- [23] Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier exposure. *International Conference on Learning Representations*, 2018. 2
- [24] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 3
- [25] Jinlei Hou, Yingying Zhang, Qiaoyong Zhong, Di Xie, Shiliang Pu, and Hong Zhou. Divide-and-assemble: Learning block-wise memory for unsupervised anomaly detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8791–8800, 2021. 6
- [26] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghighi, Robyn Ball, Katie Shpanskaya, et al. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 590–597, 2019. 2, 6
- [27] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016. 4, 8
- [28] Łukasz Kaiser, Ofir Nachum, Aurko Roy, and Samy Bengio. Learning to remember rare events. *arXiv preprint arXiv:1703.03129*, 2017. 2
- [29] Daniel S Kermany, Michael Goldbaum, Wenjia Cai, Carolina CS Valentim, Huiying Liang, Sally L Baxter, Alex McKeown, Ge Yang, Xiaokang Wu, Fangbing Yan, et al. Identifying medical diagnoses and treatable diseases by image-based deep learning. *Cell*, 172(5):1122–1131, 2018. 2, 6
- [30] Jung Uk Kim, Sungjune Park, and Yong Man Ro. Robust small-scale pedestrian detection with cued recall via memory learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3050–3059, 2021. 3
- [31] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6
- [32] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 2, 7
- [33] Ankit Kumar, Ozan Irsoy, Peter Ondruska, Mohit Iyyer, James Bradbury, Ishaan Gulrajani, Victor Zhong, Romain Paulus, and Richard Socher. Ask me anything: Dynamic memory networks for natural language processing. In *International conference on machine learning*, pages 1378–1387. PMLR, 2016. 2
- [34] Kimin Lee, Honglak Lee, Kibok Lee, and Jinwoo Shin. Training confidence-calibrated classifiers for detecting out-of-distribution samples. *International Conference on Learning Representations*, 2017. 2
- [35] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *Advances in neural information processing systems*, 31, 2018. 2
- [36] Sangho Lee, Jinyoung Sung, Youngjae Yu, and Gunhee Kim. A memory network approach for story-based temporal summarization of 360 videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1410–1419, 2018. 2
- [37] Chun-Liang Li, Kihyuk Sohn, Jinsung Yoon, and Tomas Pfister. Cutpaste: Self-supervised learning for anomaly detection and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9664–9674, 2021. 2, 6, 7
- [38] Jingyuan Li, Ning Wang, Lefei Zhang, Bo Du, and Dacheng Tao. Recurrent feature reasoning for image inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7760–7768, 2020. 4, 16

- [39] Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. *International Conference on Learning Representations*, 2017. 2
- [40] Guilin Liu, Fitsum A Reda, Kevin J Shih, Ting-Chun Wang, Andrew Tao, and Bryan Catanzaro. Image inpainting for irregular holes using partial convolutions. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 85–100, 2018. 4
- [41] Zhian Liu, Yongwei Nie, Chengjiang Long, Qing Zhang, and Guiqing Li. A hybrid video anomaly detection framework via memory-augmented flow reconstruction and flow-guided frame prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13588–13597, 2021. 2, 3
- [42] Yiwei Lu, K Mahesh Kumar, Seyed shahabeddin Nabavi, and Yang Wang. Future frame prediction using convolutional vrnn for anomaly detection. In *2019 16th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, pages 1–8. IEEE, 2019. 2
- [43] Yuhang Lu, Weijian Li, Kang Zheng, Yirui Wang, Adam P Harrison, Chihung Lin, Song Wang, Jing Xiao, Le Lu, Chang-Fu Kuo, et al. Learning to segment anatomical structures accurately from one exemplar. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 678–688. Springer, 2020. 1
- [44] Hui Lv, Chen Chen, Cui Zhen, Chunyan Xu, and Jian Yang. Learning normal dynamics in videos with meta prototype network. In *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition*, 2021. 3
- [45] Zahra Mirikharaji and Ghassan Hamarneh. Star shape prior in fully convolutional networks for skin lesion segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 737–745. Springer, 2018. 1
- [46] Sergio Naval Marimont and Giacomo Tarroni. Implicit field learning for unsupervised anomaly detection in medical images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 189–198. Springer, 2021. 2, 6
- [47] Bao Nguyen, Adam Feldman, Sarath Bethapudi, Andrew Jennings, and Chris G Willcocks. Unsupervised region-based anomaly detection in brain mri with adversarial image inpainting. In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, pages 1127–1131. IEEE, 2021. 2
- [48] Salima Omar, Asri Ngadi, and Hamid H Jebur. Machine learning techniques for anomaly detection: an overview. *International Journal of Computer Applications*, 79(2), 2013. 2
- [49] Hyunjong Park, Jongyoun Noh, and Bumsub Ham. Learning memory-guided normality for anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14372–14381, 2020. 3
- [50] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2536–2544, 2016. 4, 16
- [51] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015. 5
- [52] Tal Reiss, Niv Cohen, Liron Bergman, and Yedid Hoshen. Panda: Adapting pretrained features for anomaly detection and segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2806–2814, 2021. 2, 6, 7
- [53] Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. Deep one-class classification. In *International conference on machine learning*, pages 4393–4402. PMLR, 2018. 2
- [54] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning internal representations by error propagation. Technical report, California Univ San Diego La Jolla Inst for Cognitive Science, 1985. 3
- [55] Mohammadreza Salehi, Niousha Sadjadi, Soroosh Baselizadeh, Mohammad Hossein Rohban, and Hamid R Rabiee. Multiresolution knowledge distillation for anomaly detection. 2020. 2, 7
- [56] Mohammadreza Salehi, Niousha Sadjadi, Soroosh Baselizadeh, Mohammad H Rohban, and Hamid R Rabiee. Multiresolution knowledge distillation for anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14902–14912, 2021. 6, 8
- [57] Thomas Schlegl, Philipp Seeböck, Sebastian M Waldstein, Georg Langs, and Ursula Schmidt-Erfurth. f-anogan: Fast unsupervised anomaly detection with generative adversarial networks. *Medical Image Analysis*, 2019. 2, 3, 6, 7, 16, 18
- [58] Thomas Schlegl, Philipp Seeböck, Sebastian M Waldstein, Ursula Schmidt-Erfurth, and Georg Langs. Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In *International conference on information processing in medical imaging*, pages 146–157. Springer, 2017. 2, 3
- [59] Thomas Schlegl, Sebastian M Waldstein, Wolf-Dieter Vogl, Ursula Schmidt-Erfurth, and Georg Langs. Predicting semantic descriptions from medical images with convolutional neural networks. In *International Conference on Information Processing in Medical Imaging*, pages 437–448. Springer, 2015. 2
- [60] Bernhard Schölkopf, Robert C Williamson, Alexander J Smola, John Shawe-Taylor, John C Platt, et al. Support vector method for novelty detection. In *NIPS*, volume 12, pages 582–588. Citeseer, 1999. 2
- [61] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017. 7
- [62] Md Mahfuzur Rahman Siddiquee, Zongwei Zhou, Nima Tajbakhsh, Ruibin Feng, Michael B Gotway, Yoshua Ben-

- gio, and Jianming Liang. Learning fixed points in generative adversarial networks: From image-to-image translation to disease detection and localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 191–200, 2019. 2, 4
- [63] Desire Sidibe, Shrinivasan Sankar, Guillaume Lemaitre, Mojdeh Rastgoo, Joan Massich, Carol Y Cheung, Gavin SW Tan, Dan Milea, Ecosse Lamoureux, Tien Y Wong, et al. An anomaly detection approach for the identification of dme patients using spectral domain optical coherence tomography images. *Computer methods and programs in biomedicine*, 139:109–117, 2017. 2
- [64] Krishna Kumar Singh and Yong Jae Lee. Hide-and-seek: Forcing a network to be meticulous for weakly-supervised object and action localization. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 3544–3553. IEEE, 2017. 8
- [65] Lowell M Smoger, Clare K Fitzpatrick, Chadd W Clary, Adam J Cyr, Lorin P Maletsky, Paul J Rullkoetter, and Peter J Laz. Statistical modeling to characterize relationships between knee anatomy and kinematics. *Journal of Orthopaedic Research®*, 33(11):1620–1630, 2015. 1
- [66] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, pages 6105–6114. PMLR, 2019. 3
- [67] Youbao Tang, Yuxing Tang, Yingying Zhu, Jing Xiao, and Ronald M Summers. A disentangled generative model for disease decomposition in chest x-rays via normal image synthesis. *Medical Image Analysis*, 67:101839, 2021. 2
- [68] Yuxing Tang, Xiaosong Wang, Adam P Harrison, Le Lu, Jing Xiao, and Ronald M Summers. Attention-guided curriculum learning for weakly supervised classification and localization of thoracic diseases on chest radiographs. In *International Workshop on Machine Learning in Medical Imaging*, pages 249–258. Springer, 2018. 6, 8
- [69] Maximilian E Tschuchnig and Michael Gadermayr. Anomaly detection in medical imaging—a mini review. *arXiv preprint arXiv:2108.11986*, 2021. 2
- [70] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008. 4
- [71] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017. 5
- [72] Shenzhi Wang, Liwei Wu, Lei Cui, and Yujun Shen. Glancing at the patch: Anomaly localization with global and local feature comparison. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 254–263, 2021. 6
- [73] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M Summers. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2097–2106, 2017. 8
- [74] Qi Wei, Yinhao Ren, Rui Hou, Bibo Shi, Joseph Y Lo, and Lawrence Carin. Anomaly detection for medical images based on a one-class classification. In *Medical Imaging 2018: Computer-Aided Diagnosis*, volume 10575, page 105751M. International Society for Optics and Photonics, 2018. 6
- [75] Julia Wolleb, Robin Sandkühler, and Philippe C Cattin. Descargan: Disease-specific anomaly detection with weak supervision. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 14–24. Springer, 2020. 6
- [76] Tiange Xiang, Chaoyi Zhang, Yang Song, Siqi Liu, Hongliang Yuan, and Weidong Cai. Partial graph reasoning for neural network regularization. *arXiv preprint arXiv:2106.01805*, 2021. 5, 8
- [77] Muhammad Zaigham Zaheer, Arif Mahmood, M Haris Khan, Marcella Astrid, and Seung-Ik Lee. An anomaly detection system via moving surveillance robots with human collaboration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2595–2601, 2021. 3
- [78] Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. Reconstruction by inpainting for visual anomaly detection. *Pattern Recognition*, 112:107706, 2021. 2
- [79] Xiaolin Zhang, Yunchao Wei, Jiashi Feng, Yi Yang, and Thomas S Huang. Adversarial complementary learning for weakly supervised object localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1325–1334, 2018. 8
- [80] Bin Zhao, Li Fei-Fei, and Eric P Xing. Online detection of unusual events in videos via dynamic sparse coding. In *CVPR 2011*, pages 3313–3320. IEEE, 2011. 2
- [81] He Zhao, Yuexiang Li, Nanjun He, Kai Ma, Leyuan Fang, Huiqi Li, and Yefeng Zheng. Anomaly detection for medical images using self-supervised and translation-consistent features. *IEEE Transactions on Medical Imaging*, 2021. 2, 6, 7
- [82] Tianyi Zhao, Kai Cao, Jiawen Yao, Isabella Nogues, Le Lu, Lingyun Huang, Jing Xiao, Zhaozheng Yin, and Ling Zhang. 3d graph anatomy geometry-integrated network for pancreatic mass segmentation, diagnosis, and quantitative patient management. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13743–13752, 2021. 1
- [83] Kang Zhou, Jing Li, Yuting Xiao, Jianlong Yang, Jun Cheng, Wen Liu, Weixin Luo, Jiang Liu, and Shenghua Gao. Memorizing structure-texture correspondence for image anomaly detection. *IEEE Transactions on Neural Networks and Learning Systems*, 2021. 3
- [84] Kang Zhou, Yuting Xiao, Jianlong Yang, Jun Cheng, Wen Liu, Weixin Luo, Zaiwang Gu, Jiang Liu, and Shenghua Gao. Encoding structure-texture relation with p-net for anomaly detection in retinal images. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XX 16*, pages 360–377. Springer, 2020. 6
- [85] Zongwei Zhou, Jae Shin, Ruibin Feng, R Todd Hurst, Christopher B Kendall, and Jianming Liang. Integrating active learning and transfer learning for carotid intima-media



- thickness video interpretation. *Journal of digital imaging*, 32(2):290–299, 2019. [1](#)
- [86] Zongwei Zhou, Vatsal Sodha, Jiaxuan Pang, Michael B Gotway, and Jianming Liang. Models genesis. *Medical image analysis*, 67:101840, 2021. [2](#), [4](#), [16](#)
- [87] Zongwei Zhou, Vatsal Sodha, Md Mahfuzur Rahman Siddiquee, Ruibin Feng, Nima Tajbakhsh, Michael B Gotway, and Jianming Liang. Models genesis: Generic autodidactic models for 3d medical image analysis. In *International conference on medical image computing and computer-assisted intervention*, pages 384–393. Springer, 2019. [2](#)
- [88] Arthur Zimek and Erich Schubert. Outlier detection. In *Encyclopedia of Database Systems*. Springer, 2017. [2](#)
- [89] Bo Zong, Qi Song, Martin Renqiang Min, Wei Cheng, Cristian Lumezanu, Daeki Cho, and Haifeng Chen. Deep autoencoding gaussian mixture model for unsupervised anomaly detection. In *International conference on learning representations*, 2018. [2](#)

## A. Technical Details of SQUID

In this section, we present the implementation of SQUID for both training and inference using Python pseudocode in PyTorch style. The overall SQUID (§3.1) is presented in Alg. 1; The proposed in-painting block (§3.3) is presented in Alg. 2; The use of SQUID for anomaly discrimination (§3.4) is presented in Alg. 3; Finally, the adversarial training process of SQUID with the introduced losses (§3.5) is presented in Alg. 4. Code and dataset will be available.

---

### Algorithm 1 SQUID (§3.1)

---

```
# I: input X-ray image
# Q: Memory Queue
# E: encoder network
# T: teacher network
# S: student (main) network
# D: discriminator network
# M: decoder Memory Matrices
# n_patch: number of patches

def squid(I, Q, n_patch):
    skips = []
    level_features = []

    # divide non-overlapping patches
    x = divide_patch(I, n_patch)
    x = conv2d(x, kernel_size=3, stride=1)

    ##### feature extraction #####
    for encoder_layer in E:
        x = encoder_layer(x)
        skips.append(x)
    ##### feature extraction #####

    # stop-gradients for teacher network
    t_x = x.detach()

    # feature in-painting
    x = inpaint(x, Q)

    ##### image reconstruction #####
    combined = False
    skips = skips[::-1]
    for i in range(len(S)):
        # put patches together
        if not memory_matrix[i] and not combined:
            x = reverse_patch(x, n_patch)
            combined = True
        x = S[i](x, skips[i])
        x = memory_matrix[i](x)

    # stop-gradients for teacher network
    t_x = T[i](t_x, skips[i].detach())
    # preserve for distillation loss
    level_features.append((x, t_x))

    x = conv2d(x, kernel_size=1, stride=1)
    x = sigmoid(x)

    t_x = conv2d(t_x, kernel_size=1, stride=1)
    t_x = sigmoid(t_x)
    ##### image reconstruction #####

    return x, t_x, level_features
```

---



---

### Algorithm 2 In-painting Block (§3.3)

---

```
# F: extracted patch features
# Q: Memory Queue
# mask_prob: rho in Eq. 2

def inpaint(F, Q):
    shortcut = F

    # bottleneck
    F = conv2d(F, kernel_size=1, stride=1)

    # make 3x3 neighborhoods w/o centre patch
    F = make_windows(F, window_size=(3,3))

    N = torch.zeros_like(F)
    for i in range(len(F)): # sliding window
        # Figure 5, step a
        N[i] = Q[i].query_assemble(F[i])
        # Figure 5, step c
        Q[i].update(F[i])

    # Figure 5, step b
    F = transformer_layer(query=F, key=N, \
                           value=N) + F

    # reverse 3x3 neighborhoods
    F = reverse_windows(F, window_size=(3,3))

    # bottleneck
    F = conv2d(F, kernel_size=1, stride=1)

    # only activated for training
    if self.training:
        b, c, w, h = F.shape
        # sample binary mask based on mask_prob
        mask = torch.ones(b, 1, w, h) * mask_prob
        mask = torch.bernoulli(mask)

        # masked shortcuts
        F = F * mask + shortcut * (1 - mask)

    return F
```

---



---

### Algorithm 3 Anomaly Discrimination (§3.4)

---

```
# Q: Updated Memory Queue
# n_patch: number of patches
# early_stop: step for early stop

train_scores = []

# load normal training images
for step, I in enumerate(train_loader):
    train_scores.append(D(squid(I, Q, n_patch)))
    # early stopping to save time
    if step == early_stop:
        break

test_scores = []
test_labels = []

# load normal/abnormal testing images
for step, (I, y) in enumerate(test_loader):
    test_scores.append(D(squid(I, Q, n_patch)))
    test_labels.append(y)

test_scores -= mean(train_scores)
test_scores /= std(train_scores)

results = calc_metrics(test_scores, test_labels)
```

---

---

**Algorithm 4** Training SQUID (§3.5)

---

```

# CE: cross entropy
# MSE: mean square error
# critic: generator training interval
# n_patch: number of patches
# n_pattern: number of patterns per segment

# initialize Memory Queue
Q = random_initialize(n_patch, n_pattern)

# load normal images only for training
for step, I in enumerate(loader):
    ##### optimize discriminator #####
    real = D(I)
    fake = D(squid(I, Q, n_patch).detach())
    l_real = CE(I, torch.ones_like(real))
    l_fake = CE(I, torch.zeros_like(fake))

    # w: loss weight
    l_dis = w_dis * (l_real + l_fake)
    l_dis.backward()
    ##### optimize discriminator #####

    if step % critic == 0:
        x, t_x, level_features = squid(I, Q, n_patch)

        ##### loss calculations #####
        # w: the loss weights
        l_s = w_s * MSE(I, x)
        l_t = w_t * MSE(I, t_x)
        l_dist = w_dist * MSE(level_features)

        fake = D(x)
        l_gen = w_gen * \
            CE(fake, torch.ones_like(fake))

        loss = l_s + l_t + l_dist + l_gen
        loss.backward()
        ##### loss calculations #####

```

---

## B. Architectures of SQUID

Our SQUID consists of an encoder, a student (main) generator, a teacher generator, and a discriminator. All of the network architectures are built with plain convolution, batch normalization, and ReLU activation layers only. The architecture details of the encoder are shown in Table 3. For an input X-ray image (sized of  $128 \times 128$ ), we first divide it into  $2 \times 2$  non-overlapping patches (sized of  $64 \times 64$ ). The encoder then extracts the patch features.

As mentioned in §3.1, the student and teacher generators were constructed identically. The only difference is that additional Memory Matrices are placed in the student generator. The architecture details of the student generator are shown in Table 4. Skip connections from the encoder are only enabled at such levels that Memory Matrices are used. After the last Memory Matrix, the non-overlapping patches are put back as a whole for further reconstruction.

As shown in Table 5, the discriminator was constructed in a more lightweight style. Note that the images are discriminated at their full resolution (*i.e.*  $128 \times 128$ ) rather than in patches.

Level	#Channels	Resolution
Input	1	$(2 \times 2) \times (64 \times 64)$
1	32	$(2 \times 2) \times (32 \times 32)$
2	64	$(2 \times 2) \times (16 \times 16)$
3	128	$(2 \times 2) \times (8 \times 8)$
4	256	$(2 \times 2) \times (4 \times 4)$

Table 3. Encoder structure in SQUID.

Level	#Channels	w/ S&M	Resolution
4	256	✓	$(2 \times 2) \times (4 \times 4)$
3	128	✓	$(2 \times 2) \times (8 \times 8)$
2	64		$32 \times 32$
1	32		$64 \times 64$
Output	1		$128 \times 128$

Table 4. Student and teacher generator structures in SQUID. S&amp;M denotes the usage of skip connections and Memory Matrix. There is no Memory Matrix placed in the teacher generator.

Level	#Channels	Resolution
Input	1	$128 \times 128$
1	16	$64 \times 64$
2	32	$32 \times 32$
3	64	$16 \times 16$
4	128	$8 \times 8$
5	128	$4 \times 4$
Output	1	$1 \times 1$

Table 5. Discriminator structure in SQUID.

## C. Additional Results

### C.1. Extensive Ablation Studies

In this section, we ablate three more designs in SQUID to fully validate their necessity and effectiveness.

**(1) Convolutional vs. Transformer Layers:** In our proposed in-painting block, a transformer layer is used to aggregate the encoder extracted patch features and the Memory Queue augmented “normal” features. However, one may wonder if a simple convolution layer can also suffice. We conducted experiments by replacing the transformer layer with a convolutional layer while preserving other structures.

**(2) Soft vs. Hard Masked Shortcuts:** In our proposed masked shortcut, skipped and in-painted features are aggregated using a binary gating mask. The intuitive question is whether such “hard” gating is necessary and a weighted “soft” addition can also achieve comparable results. To this end, instead of following Eq. 2, we conducted experiments by aggregating the patch features  $\mathcal{F}$  through:

$$\mathcal{F}' = (1 - \rho) \cdot \mathcal{F} + \rho \cdot \text{inpaint}(\mathcal{F}), \quad (3)$$

where  $\rho$  was set to 95%, following the best setting adopted in SQUID.

**(3) Pixel-level vs. Feature-level In-painting:** As discussed in §3.3, raw images usually contain larger noise

Method	AUC (%)	Acc (%)	F1 (%)
Convolution Layers	76.9±3.3	74.2±3.3	80.7±2.7
Transformer Layers ( $\Delta$ )	$\uparrow 10.7$	$\uparrow 6.1$	$\uparrow 4.0$
Soft Masked Shortcut	79.7±3.4	76.1±2.7	80.7±2.3
Hard Masked Shortcut ( $\Delta$ )	$\uparrow 7.9$	$\uparrow 4.2$	$\uparrow 4.0$
Pixel-level In-painting	79.1±0.4	74.4±1.6	81.3±0.9
Feature-level In-painting ( $\Delta$ )	$\uparrow 8.5$	$\uparrow 5.9$	$\uparrow 3.4$
Full SQUID	<b>87.6±1.5</b>	<b>80.3±1.3</b>	<b>84.7±0.8</b>

Table 6. The extensive results indicate that all proposed techniques in SQUID are essential for a high overall performance.

#Patches	AUC (%)	Acc (%)	F1 (%)	Sen. (%)	Spec. (%)
1×1	77.6±0.5	75.5±0.5	82.5±0.6	<b>92.6±0.6</b>	47.0±0.7
2×2	<b>87.6±1.5</b>	<b>80.3±1.3</b>	<b>84.7±0.8</b>	87.8±3.3	<b>70.2±3.6</b>
3×3	82.9±2.4	77.5±1.3	82.8±0.9	86.0±1.1	61.5±1.5
4×4	83.5±1.3	79.0±1.5	83.8±0.7	89.5±2.0	56.0±2.9
8×8	82.0±0.9	77.6±1.9	83.0±1.8	87.9±2.8	60.4±3.4

Table 7. Ablation on the number of patch division.

Top $K$	AUC (%)	Acc (%)	F1 (%)	Sen. (%)	Spec. (%)
1	72.6±2.0	71.7±2.3	78.6±2.7	83.7±7.0	51.9±5.8
5	<b>87.6±1.5</b>	<b>80.3±1.3</b>	<b>84.7±0.8</b>	87.8±3.3	<b>70.2±3.6</b>
10	82.8±4.0	77.8±2.5	83.0±1.9	86.2±2.0	63.8±3.5
50	82.4±1.9	76.9±2.2	82.8±1.3	<b>89.1±4.6</b>	56.7±11.3
100	82.0±5.1	76.7±3.1	82.5±2.3	88.2±4.6	57.4±9.5

Table 8. Ablation on the top  $K$  value in the Gumbel Shrinkage.

and artifacts than features, so we proposed to achieve the in-painting at the feature level rather than at the image level [38, 50, 86]. To validate our claim, we have conducted experiments on carrying out the in-painting at the pixel level. Instead of using a transformer layer to in-paint the extracted patch features, we randomly zeroed out parts of the input patches with 25% probability and let SQUID in-paint the distorted input images. All other settings and objective functions remain unchanged.

**Summary:** The results of the above three additional ablation experiments are presented in Table 6. Without using the transformer layer, masked shortcut, and feature-level in-painting as proposed, the AUC, Acc, and F1 scores decreased by at least 8%, 4%, and 3%, respectively, compared with the full SQUID setting.

## D. Tabular Results

In this section, we present the tabular results of the ablation studies introduced in Figure 9 and offer additional metrics (*i.e.* Accuracy, F1, Sensitivity, and Specificity). The results containing the number of patch division, value of the topk in Gumbel Shrinkage, and mask probability in masked shortcuts are shown in Table 7, Table 8, and Table 9.

## E. Creating DigitAnatomy

The pseudocode of creating our new benchmark dataset (DigitAnatomy in §4.1) is provided in Alg. 5. In practice,

$\rho$	AUC (%)	Acc (%)	F1 (%)	Sen. (%)	Spec. (%)
100%	82.5±1.3	76.4±0.8	82.3±1.1	87.4±0.7	60.8±2.7
99%	82.0±0.8	75.6±0.6	80.8±1.1	82.1±0.6	65.0±1.6
95%	<b>87.6±1.5</b>	<b>80.3±1.3</b>	<b>84.7±0.8</b>	87.8±3.3	<b>70.2±3.6</b>
90%	82.9±0.6	79.0±0.9	83.7±1.4	87.3±1.1	60.2±2.9
85%	79.0±2.3	76.9±1.2	82.8±1.4	<b>91.4±2.6</b>	50.0±3.6
80%	79.3±2.2	75.7±1.7	82.0±1.4	88.5±1.7	54.4±4.9
70%	77.9±2.8	74.4±3.0	80.8±1.7	86.0±2.1	55.3±7.0
60%	79.3±3.5	75.8±2.1	81.6±0.7	85.9±3.3	58.9±8.0
50%	76.0±5.3	74.7±2.2	81.1±0.5	86.7±5.1	54.8±9.4

Table 9. Ablation on mask probability  $\rho$  in the masked shortcuts.

### Algorithm 5 Creating DigitAnatomy

```

# a function to pick random digit instances
def pick_random(class, single_digits):
    # random pick an image with size: [28, 28]
    pick_digit = random.choice(single_digits[class])
    return pick_digit

# load MNIST digits with shape: [10, 1000, 28, 28]
single_digits = load_MNIST()

# all possible conditions
conditions = ['normal', 'missing', \
             'misorder', 'flipped', 'novel']

output = torch.zeros(3, 28, 3, 28)

# loop over digit 1-9 in order
for idx in range(1, 10):

    # randomly pick a condition
    condition = random.choice(conditions)

    if condition == 'normal':
        digit = pick_random(idx, single_digits)
    # anatomy of missing digit
    elif condition == 'missing':
        digit = torch.zeros(28, 28)
    # anatomy of disorder digit
    elif condition == 'misorder':
        ridx = random.randint(1, 10)
        digit = pick_random(ridx, single_digits)
    # anatomy of flipped digit
    elif condition == 'flipped':
        digit = pick_random(idx, single_digits)
        digit = digit[:, ::-1, ::-1]
    # anatomy of novel digit
    elif condition == 'novel':
        digit = pick_random(0, single_digits)

    output[idx // 3, :, idx % 3, :] = digit

# combine all patches together
output = output.view(28 * 3, 28 * 3)

```

we have implemented the algorithm into an off-the-shelf data loader that can be amended to many other different datasets (*e.g.* SVHN, CIFAR, ImageNet).

## F. Visualization Results

### F.1. Visualizations on DigitAnatomy

More reconstruction results of SQUID and the compared methods [1, 13, 57] are shown in Figure 10. Our obser-



variations from these additional results are aligned with the ones discussed in §5.1. SQUID can capture *every* appearing anomaly (highlighted in red) in the images and augment them back to the normal closest forms. On the contrary, although MemAE restores the normal digits the best, it is limited in detecting a few anomaly types (*e.g.* misordered and missing digits). Ganomaly is not able to perfectly recover the normal digits and also cannot generate meaningful reconstructions on the abnormal ones. f-AnoGAN, on the other hand, memorizes and generates an exemplary normal pattern that fails to respond to different inputs.

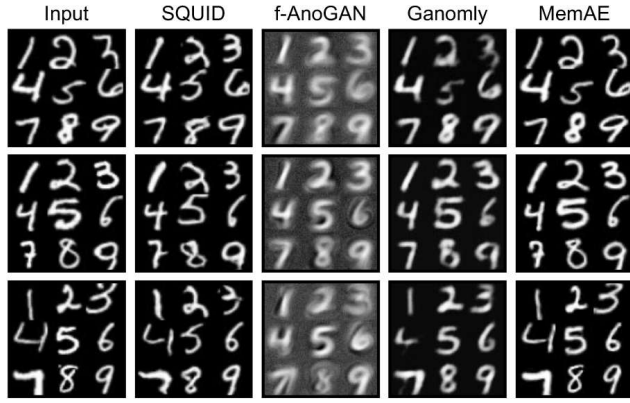
## F.2. Visualizations on Chest X-rays

Figure 11 and Figure 12 show more reconstruction results of our SQUID on the ZhangLab Chest X-ray and Stanford CheXpert datasets. We observed that our method is capable of augmenting the input images to similar “normal” outputs and assigns larger anomaly scores to abnormal cases.

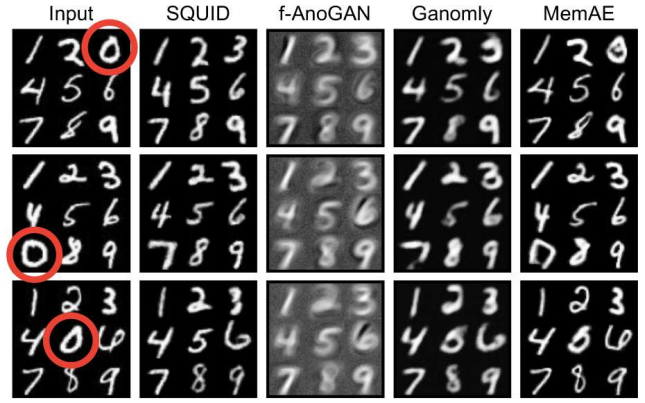
When inputting normal images, SQUID will try to reconstruct the inputs as well as possible. Due to the usage of memory modules, our framework could hardly degenerate to function as an identity mapping from inputs to outputs. Therefore, the reconstruction of normal inputs cannot perfectly recover every single detail.

When inputting abnormal images, SQUID will make larger impacts by combining previously seen normal features together into such abnormal ones. Since the generator is not trained on such hybrid features, the reconstruction results could demonstrate more obvious artifacts and blurs.

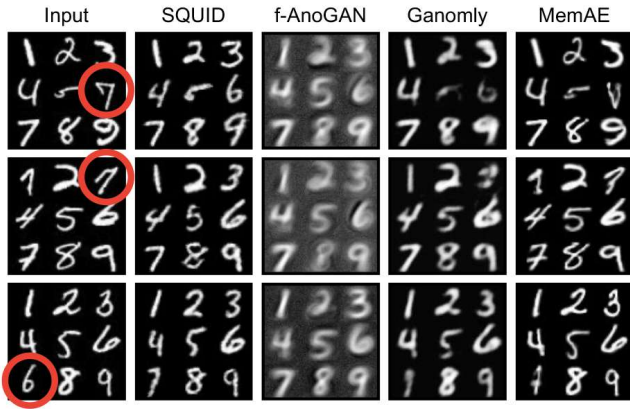
After our framework converges, the optimized discriminator can perceptually capture such inconsistencies between reconstructed normal and abnormal images and achieve anomaly detection.



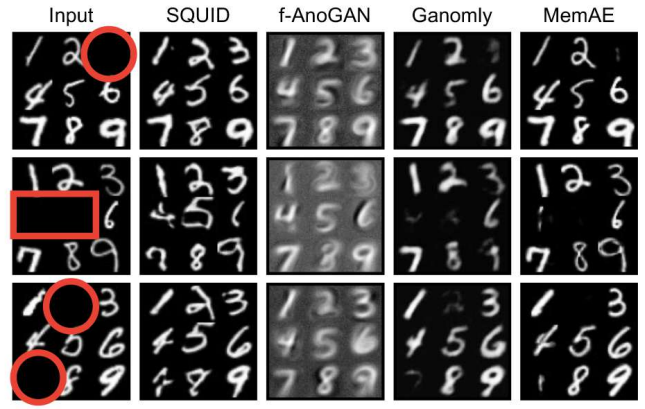
(a) Normal



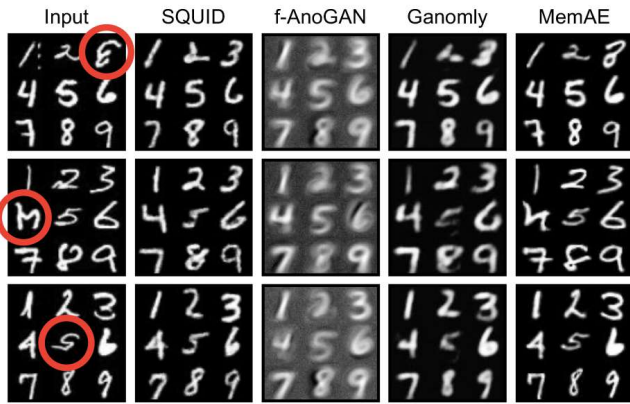
(b) Abnormal (novel digit)



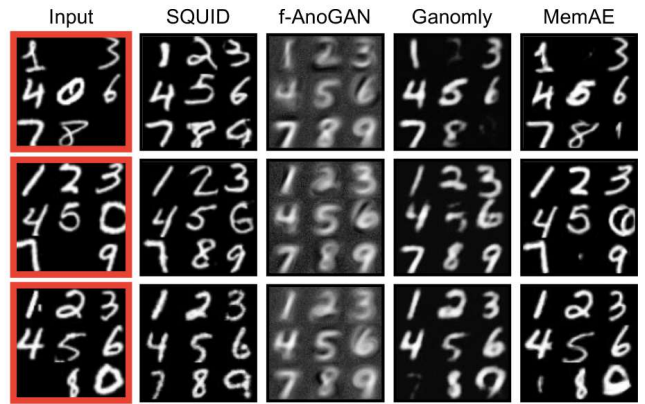
(c) Abnormal (misordered)



(d) Abnormal (missing digit)



(e) Abnormal (flipped digit)



(f) Abnormal (mixture)

Figure 10. Comparisons of reconstruction results on DigitAnatomy of our SQUID, f-AnoGAN [57], Ganomly [1], and MemAE [13]. Anomalies are highlighted in red.

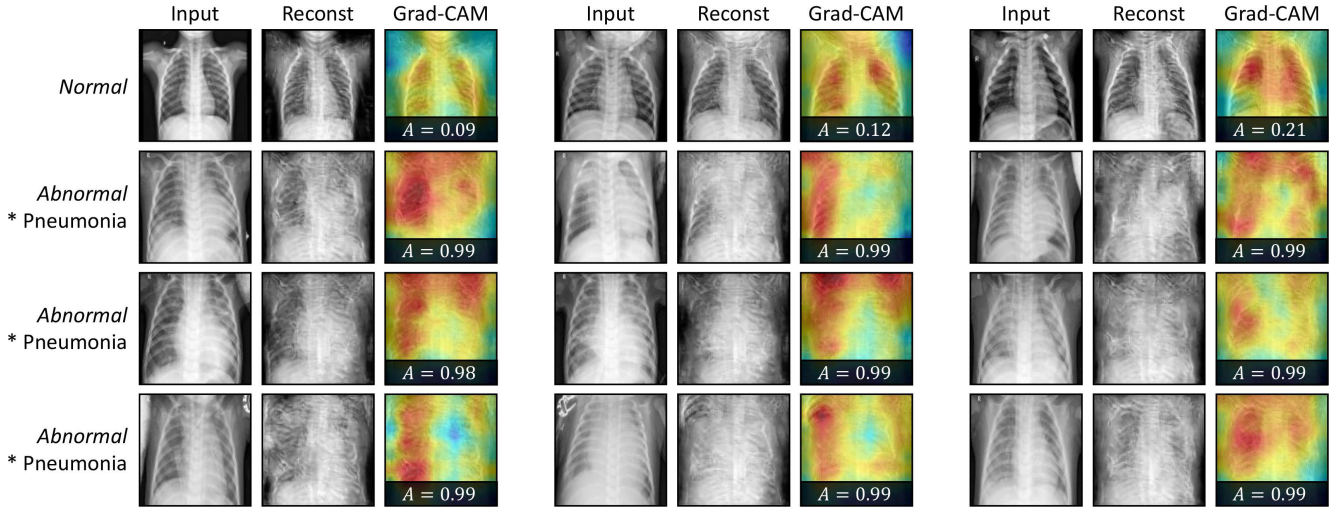


Figure 11. Reconstruction results of SQUID on the ZhangLab Chest X-ray dataset. The corresponding Grad-CAM heatmaps along with anomaly scores are shown as well.

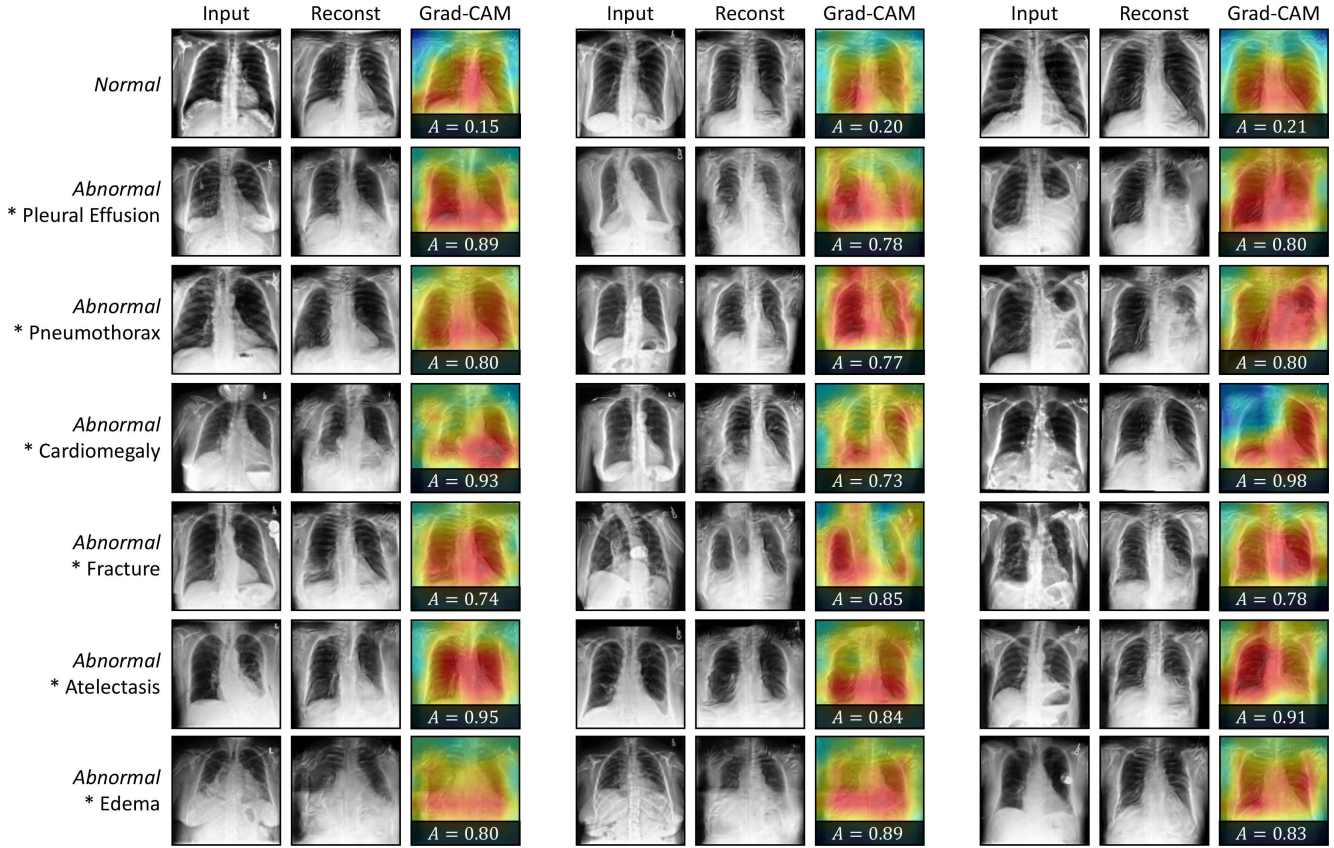


Figure 12. Reconstruction results of SQUID on the Stanford CheXpert dataset. Different disease types are separated in different rows. The corresponding Grad-CAM heatmaps along with anomaly scores are shown as well.