

DiAD: A Diffusion-based Framework for Multi-class Anomaly Detection

Haoyang He^{1*}, Jiangning Zhang^{2*}, Hongxu Chen¹, Xuhai Chen¹, Zhishan Li¹,
Xu Chen², Yabiao Wang², Chengjie Wang², Lei Xie^{1†}

¹Zhejiang University ²Youtu Lab, Tencent

Abstract

Reconstruction-based approaches have achieved remarkable outcomes in anomaly detection. The exceptional image reconstruction capabilities of recently popular diffusion models have sparked research efforts to utilize them for enhanced reconstruction of anomalous images. Nonetheless, these methods might face challenges related to the preservation of image categories and pixel-wise structural integrity in the more practical multi-class setting. To solve the above problems, we propose a *Diffusion-based Anomaly Detection (DiAD)* framework for multi-class anomaly detection, which consists of a pixel-space autoencoder, a latent-space *Semantic-Guided (SG)* network with a connection to the stable diffusion’s denoising network, and a feature-space pre-trained extractor. Firstly, The SG network is proposed for reconstructing anomalous regions while preserving the original image’s semantic information. Secondly, we introduce *Spatial-aware Feature Fusion (SFF)* block to maximize reconstruction accuracy when dealing with extensively reconstructed areas. Thirdly, the input and reconstructed images are processed by a pre-trained feature extractor to generate anomaly maps based on features extracted at different scales. Experiments on MVTEC-AD and VisA datasets demonstrate the effectiveness of our approach which surpasses the state-of-the-art methods, *e.g.*, achieving 96.8/52.6 and 97.2/99.0 (AUROC/AP) for localization and detection respectively on multi-class MVTEC-AD dataset. Code will be available at <https://lewandofskoe.github.io/projects/diad>.

Introduction

Anomaly detection is a crucial task in computer vision and industrial applications (Tao et al. 2022; Salehi et al. 2022; Liu et al. 2023), which goal of visual anomaly detection is to determine anomalous images and locate the regions of anomaly accurately. Existing anomaly detection models (Liznerski et al. 2021; Yi and Yoon 2020; Yu et al. 2021) mostly correspond to one class, which requires a large amount of storage space and training time as the number of classes increases. Therefore, there is an urgent need for an unsupervised multi-class anomaly detection model that is robust and stable.

The current mainstream unsupervised anomaly detection methods can be divided into three categories: synthesizing-

*Equal contribution.

†Corresponding author.

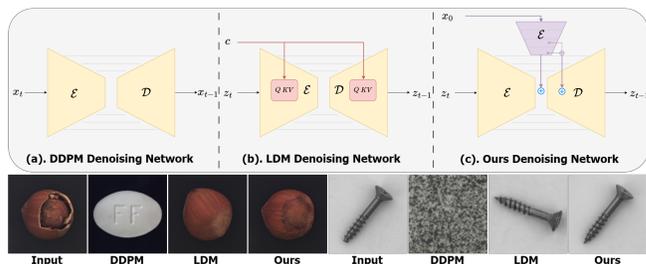


Figure 1: A analysis of different diffusion models for multi-class anomaly detection. The image above shows various denoising network architectures, while the images below demonstrate the results reconstructed by different methods for the same input image. *a)* DDPM suffers from categorical errors. *b)* LDM exhibits semantic errors. *c)* Our approach effectively reconstructs the anomalous regions while preserving the semantic information of the original image.

based (Zavrtanik, Kristan, and Skočaj 2021a; Li et al. 2021), embedding-based (Defard et al. 2021; Roth et al. 2022; Xie et al. 2023) and reconstruction-based (Liu et al. 2022; Liang et al. 2023) methods. The central concept of the reconstruction-based method is that during the training phase, the model only learns from normal images. During the testing phase, the model reconstructs abnormal images into normal ones using the trained model. Therefore, by comparing the reconstructed image with the input image, we can determine the location of anomalies. Traditional reconstruction-based methods, including AEs (Zavrtanik, Kristan, and Skočaj 2021b), VAEs (Kingma and Welling 2022), and GANs (Liang et al. 2023; Yan et al. 2021) can learn the distribution of normal samples and reconstruct abnormal regions during the testing phase. However, these models have limited reconstruction capabilities and cannot reconstruct complicated textures and objects well, especially large-scale defects or disappearances as shown in Figure 1. Hence, models with stronger reconstruction capability are required to effectively tackle multi-class anomaly detection.

Recently, the diffusion models (Ho, Jain, and Abbeel 2020; Rombach et al. 2022; Zhang and Agrawala 2023) have demonstrated their powerful image-generation capability. However, directly using current mainstream diffusion

models cannot effectively address multi-class anomaly detection problems. 1) For the Denoising Diffusion Probabilistic Model (DDPM) (Ho, Jain, and Abbeel 2020) in Fig. 1-(a), when performing the multi-class setting, this method may encounter issues with misclassifying generated image categories. The reason is that after adding T timesteps noise to the input image, the image has lost its original class information. During inference, denoising is performed based on this Gaussian noise-like distribution, which may generate samples belonging to different categories. 2) Latent Diffusion Model (LDM) (Rombach et al. 2022) has an embedder as a class condition as shown in Fig. 1-(b), which does not exist the problem of misclassification found in DDPM. However, LDM still cannot address the issue of semantic loss in generated images. LDM is unable to simultaneously preserve the semantic information of the input image while reconstructing the anomalous regions. For example, they may fail to maintain direction consistency with the input image in terms of objects like screws and hazelnuts, as well as exhibit substantial differences from the original image in terms of texture class images.

To address the aforementioned problems, we propose a diffusion-based framework, DiAD, for multi-class anomaly detection and localization, illustrated in Fig. 2, which comprises three components: a pixel space autoencoder, a latent space denoising network and a feature space ImageNet pre-trained model. To effectively maintain consistent semantic information with the original image while reconstructing the location of anomalous regions, we propose the Semantic-Guided (SG) network with a connection to the Stable Diffusion (SD) denoising network in LDM. To further enhance the capability of preserving fine details in the original image and reconstructing large defects, we propose the Spatial-aware Feature Fusion (SFF) block to integrate features at different scales. Finally, the reconstructed and input images are passed through a pre-trained model to extract features at different scales and compute anomaly scores. We summarize our contributions as follows:

- We propose a novel diffusion-based framework DiAD for multi-class anomaly detection, which firstly tackles the problem of existing denoising networks of diffusion-based methods failing to correctly reconstruct anomalies.
- We construct an SG network connecting to the SD denoising network to maintain consistent semantic information and reconstruct the anomalies.
- We propose an SFF block to integrate features from different scales to further improve the anomaly reconstruction ability.
- Abundant experiments demonstrate the sufficient superiority of DiAD over SOTA methods, *e.g.*, we surpass the multi-class anomaly detection diffusion-based method by 20.6 \uparrow / 11.7 \uparrow in pixel/image AUROC and non-diffusion method by 9.2 \uparrow in pixel-AP and 0.7 \uparrow in image-AUROC on MVTEC-AD dataset.

Related work

Diffusion model. The diffusion model has gained widespread attention and research interest since its remark-

able reconstruction ability. It has demonstrated excellent performance in various applications such as image generation (Zhang and Agrawala 2023), video generation (Ho et al. 2022), object detection (Chen et al. 2022), image segmentation (Amit et al. 2022) and etc. LDM (Rombach et al. 2022) introduces conditions through cross-attention to control generation. However, it fails to accurately reconstruct images that contain the original semantic information.

Anomaly detection. AD contains a variety of different settings, *e.g.*, open-set (Ding, Pang, and Shen 2022), noisy learning (Tan et al. 2021; Yoon et al. 2022), zero-/few-shot (Huang et al. 2022; Jeong et al. 2023; Cao et al. 2023; Chen, Han, and Zhang 2023; Chen et al. 2023b; Zhang et al. 2023b), 3D AD (Wang et al. 2023; Chen et al. 2023a), *etc.* This paper studies general unsupervised anomaly detection, which can primarily be categorized into three major methodologies:

1) Synthesizing-based methods synthesize anomalies on normal image samples. During the training phase, both normal images and synthetically generated abnormal images are input into the network for training, which aids in anomaly detection and localization. DRAEM (Zavrtanik, Kristan, and Skočaj 2021a) consists of an end-to-end network composed of a reconstruction network and a discriminative sub-network, which synthesizes and generates just-out-distribution phenomena. However, due to the diversity and unpredictability of anomalies in real-world scenarios, it is impossible to synthesize all types of anomalies.

2) Embedding-based methods encode the original image’s three-dimensional information into a multidimensional feature space (Roth et al. 2022; Cao et al. 2022; Gu et al. 2023). Most methods employ networks (He et al. 2016; Tan and Le 2019; Zhang et al. 2022, 2023c; Wu et al. 2023) pre-trained on ImageNet (Deng et al. 2009) for feature extraction. RD4AD (Deng and Li 2022) utilizes a WideResNet50 (Zagoruyko and Komodakis 2016) as the teacher model for feature extraction and employs a structurally identical network in reverse as the student model, computing the cosine similarity of corresponding features as anomaly scores. However, due to significant differences between industrial images and the data distribution in ImageNet, the extracted features might not be suitable for industrial anomaly detection purposes.

3) Reconstruction-based methods aim to train a model on a dataset without anomalies. The model learns to identify patterns and characteristics in the normal data. OCR-GAN (Liang et al. 2023) decouples images into different frequencies and uses GAN for reconstruction. EdgRec (Liu et al. 2022) achieves good reconstruction results by first synthesizing anomalies and then extracting grayscale edge information from images, which is ultimately input into a reconstruction network. However, there are certain limitations in the reconstruction of large-area anomalies. Moreover, the accuracy of anomaly localization is also not sufficient.

Recently, some studies have applied diffusion models to anomaly detection. AnoDDPM (Wyatt et al. 2022) is the first approach to employ a diffusion model for medical anomaly detection. DiffusionAD (Zhang et al. 2023a) utilizes an

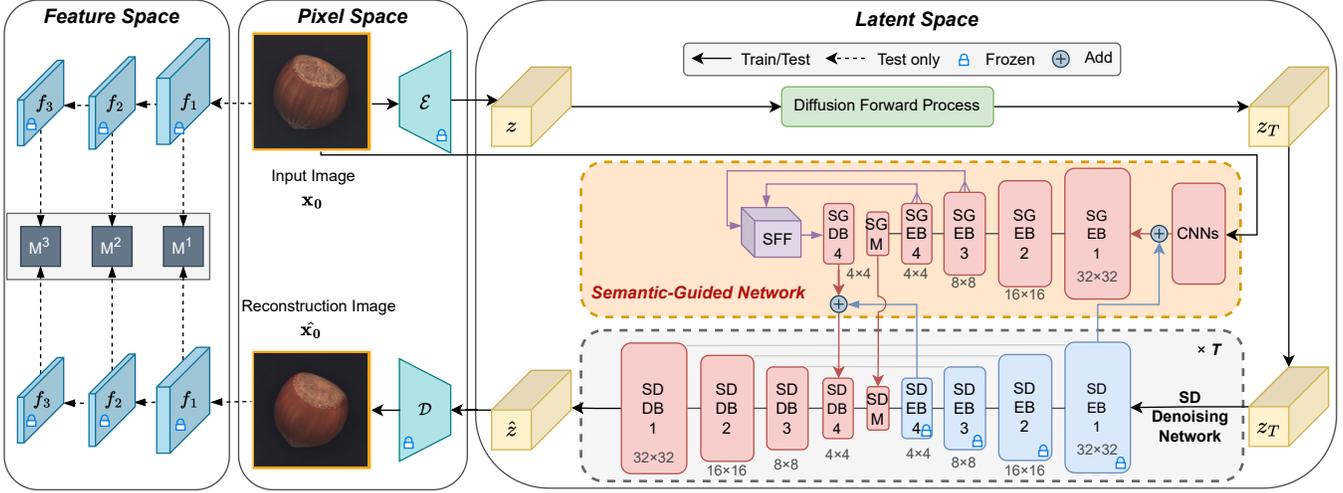


Figure 2: **Framework of the proposed DiAD that contains three parts:** 1) a pixel-space autoencoder $\{\mathcal{E}, \mathcal{D}\}$; 2) a latent-space Semantic-Guided (SG) network with a connection to Stable Diffusion (SD) denoising network; and 3) a feature-space pre-trained feature extractor Ψ . During training, the input x_0 and the latent variable z_T are inputted into the SG network and the SD denoising network, respectively. The output noise and input noise are calculated for MSE loss and gradient optimization is computed. During testing, x_0 and the reconstructed image \hat{x}_0 are inputted into the same pre-trained feature extraction network to obtain feature maps $\{f_1, f_2, f_3\}$ of different scales, and their anomaly scores \mathcal{S} are calculated.

anomaly synthetic strategy to generate anomalous samples and labels, along with two sub-networks dedicated to the tasks of denoising and segmentation. DDAD (Mousakhan, Brox, and Tayyub 2023) employs a score-based pre-trained diffusion model to generate normal samples while fine-tuning the pre-trained feature extractor to achieve domain transfer. However, these approaches only add limited steps of noise and perform few denoising steps, which makes them unable to reconstruct large-scale defects.

To overcome the aforementioned problems, We propose a diffusion-based framework DiAD for multi-class anomaly detection, which firstly tackles the problem of existing diffusion-based methods failing to correctly reconstruct anomalies.

Preliminaries

Denoising Diffusion Probabilistic Model. Denoising Diffusion Probabilistic Model (DDPM) consists of two processes: the forward diffusion process and the reverse denoising process. During the forward process, a noisy sample x_t is generated using a Markov chain that incrementally adds Gaussian-distributed noise to an initial data sample x_0 . The forward diffusion process can be characterized as follows:

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t, \epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (1)$$

where $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i = \prod_{i=1}^t (1 - \beta_i)$ and β_i represents the noise schedule used to regulate the quantity of noise added at each timestep.

In the reverse denoising process, x_T is first sampled from equation 1 and x_{t-1} is reconstructed from x_t and the model prediction $\epsilon_\theta(x_t, t)$ with the formulation:

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right) + \sigma_t z, \quad (2)$$

where $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, σ_t is a fixed constant related to the variance schedule, $\epsilon_\theta(x_t, t)$ is a U-Net (Ronneberger, Fischer, and Brox 2015) network to predict the distribution and θ is the learnable parameter which could be optimized as:

$$\min_{\theta} \mathbb{E}_{x_0 \sim q(x_0), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t} \|\epsilon - \epsilon_\theta(x_t, t)\|_2^2. \quad (3)$$

Latent Diffusion Model. Latent Diffusion Model (LDM) focuses on the low-dimensional latent space with conditioning mechanisms. LDM consists of a pre-trained autoencoder model and a denoising U-Net-like attention-based network. The network compresses images using an encoder, conducts diffusion and denoising operations in the latent representation space, and subsequently reconstructs the images back to the original pixel space using a decoder. The training optimization objective is:

$$\mathcal{L}_{LDM} = \mathbb{E}_{z_0, t, c, \epsilon \sim \mathcal{N}(0, 1)} \left[\|\epsilon - \epsilon_\theta(z_t, t, c)\|_2^2 \right], \quad (4)$$

where c represents the conditioning mechanisms which can consist of multimodal types such as text or image, connected to the model through a cross-attention mechanism. z_t represents the latent space variable,

Method

The proposed pipeline DiAD is shown in Fig. 2. First, the pre-trained encoder downsamples the input image into a latent-space representation. Then, noise is added to the latent representation, followed by the denoising process using an SD denoising network with a connection to the SG network. The denoising process is repeated for the same timesteps as the diffusion process. Finally, the reconstructed latent representation is restored to the original image level using the

pre-trained decoder. In terms of anomaly detection and localization, the input and reconstructed images are fed into the same pre-trained model to extract features at different scales and calculate the differences between these features.

Semantic-Guided Network

As discussed earlier, DDPM and LDM each have specific problems when addressing multi-class anomaly detection tasks. In response to these issues and the multi-class task itself, we propose an SG network to address the problem of LDM’s inability to effectively reconstruct anomalies and preserve the semantic information of the input image.

Given an input image $x_0 \in \mathbb{R}^{3 \times H \times W}$ in pixel space, the pre-trained encoder \mathcal{E} encodes x_0 into a latent space representation $z = \mathcal{E}(x_0)$ where $z \in \mathbb{R}^{c \times h \times w}$. Similar to Eq. 1 where the original pixel space x is replaced by latent representation z , the forward diffusion process now can be characterized as follows:

$$z_t = \sqrt{\alpha_t} z_0 + \sqrt{1 - \alpha_t} \epsilon_t, \epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (5)$$

The perturbed representation z_T and input x_0 are simultaneously fed into the SD denoising network and SG network, respectively. After T steps of the reverse denoising process, the final variable \hat{z} is restored to the reconstructed image \hat{x}_0 from the pre-trained decoder \mathcal{D} giving $\hat{x}_0 = \mathcal{D}(\hat{z})$. The training objective of DiAD is:

$$\mathcal{L}_{DiAD} = \mathbb{E}_{z_0, t, c_i, \epsilon \sim \mathcal{N}(0,1)} \left[\|\epsilon - \epsilon_\theta(z_t, t, c_i)\|_2^2 \right]. \quad (6)$$

The denoising network consists of a pre-trained SD denoising network and an SG network that replicates the SD parameters for initiation as shown in Fig. 2. The pre-trained SD denoising network comprises four encoder blocks, one middle block and four decoder blocks. Here, ‘block’ means a frequently utilized unit in the construction of the neural network layer, *e.g.*, ‘resnet’ block, transformer block, multi-head cross attention block, *etc.*

The input image $x_0 \in \mathbb{R}^{3 \times H \times W}$ is transformed into $x \in \mathbb{R}^{d \times h \times w}$ by a set of ‘conv-silu’ layers \mathcal{C} in SG network in order to keep the same dimension with the latent representations in SD Encoder Block 1 \mathcal{E}_{SD1} . Then, the result of the summation of x and z are input into the SG Encoder Blocks (SGBEs). After continuous downsampling by the encoder \mathcal{E}_{SG} , the results are finally added to the output of the SD middle block \mathcal{M}_{SD} after its completion in the middle block \mathcal{M}_{SG} . Additionally, to address multi-class tasks of different scenarios and categories, the results of the SG Decoder Blocks (SGDBs) \mathcal{D}_{SG} are also added to the results of the SD decoder \mathcal{D}_{SD} with an SFF block combined which will be particularly explained in the next section. The output \mathcal{G} of the denoising network is characterized as:

$$\mathcal{G} = \mathcal{D}_{SD}(\mathcal{M}_{SD}(\mathcal{E}_{SD}(z_t)) + \mathcal{M}_{SG}(\mathcal{E}_{SD}(z + \mathcal{C}(x_0)))) + \mathcal{D}_{SG_j}(\mathcal{M}_{SG}(\mathcal{E}_{SD}(z + \mathcal{C}(x_0)))), \quad (7)$$

where z represents the latent representation with noise perturbed, x_0 represents the input image, $\mathcal{C}(\cdot)$ represents a set of ‘conv-silu’ layers in SG network, $\mathcal{E}_{SD}(\cdot)$ represents all the SD encoder blocks (SDEBs), $\mathcal{E}_{SG}(\cdot)$ represents all the

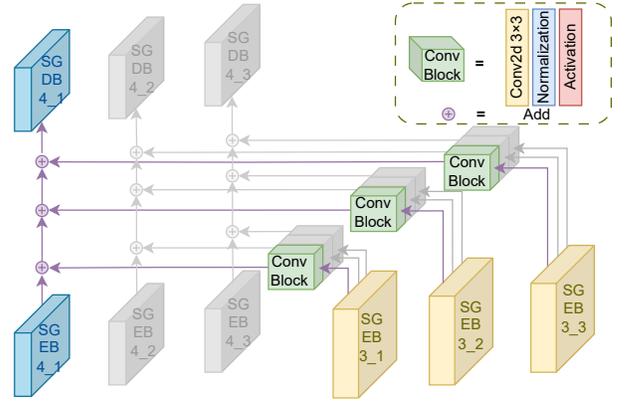


Figure 3: **Schematic diagram of SFF block.** Each layer in SGDB4 is obtained by adding the corresponding SGEB4 to every SGEB3 with Conv Block performed.

SGBEs, $\mathcal{M}_{SG}(\cdot)$ and $\mathcal{M}_{SD}(\cdot)$ represent SG and SD middle blocks respectively, $\mathcal{D}_{SD}(\cdot)$ represent all the SDBs and $\mathcal{D}_{SG_j}(\cdot)$ represents SGDBs for j -th blocks.

Spatial-aware Feature Fusion Block

When adding several layers of decoder blocks from SGBEs to SDBs during the experiment as shown in Table 7, we found it to be challenging to solve the multi-class anomaly detection. This is because the dataset contains various types, such as objects and textures. For texture-related cases, the anomalies are generally smaller, so it is necessary to preserve their original textures. On the other hand, the defects often cover larger areas for object-related cases, requiring stronger reconstruction capabilities. Therefore, it is extremely challenging to simultaneously preserve the normal information of the original samples and reconstruct the abnormal locations in different scenarios.

Hence, we proposed a Spatial-aware Feature Fusion (SFF) block with the aim of integrating high-scale semantic information into the low-scale. This ultimately enables the model to both preserve the information of the original normal samples and reconstruct large-scale abnormal regions. The structure of the SFF block is shown in Fig. 3. Each SGBEs consists of three sub-layers. Therefore, the SFF block integrates the features of each layer in SGEB3 into each layer in SGEB4 and adds the fused features to the original features. The final output of each layer of the SGEB4 is:

$$Q_i = P_i + \sum_{j=1}^J \mathcal{F}(\mathcal{H}_j), \quad (8)$$

where P_i represents the low-scale output features of the i -th layer of SGEB4, Q_i represents the final low-scale output features of the i -th layer of SGDB4, \mathcal{H}_j represents the high-scale output features of the j -th layer of SGEB3, $J = 3$ indicates three layers of SGEB3 used in the experiment and $\mathcal{F}(\cdot)$ represent a basic convolutional block which consists of a 3x3 convolution layer followed by a normalization layer and an activation layers.

Category		Non-Diffusion Method					Diffusion-based Method		
		PaDiM	MKD	DRAEM	RD4AD	UniAD	DDPM	LDM	Ours
Objects	Bottle	97.9/-	98.7/-	97.5/99.2/96.1	99.6/99.9/98.4	99.7/100./100.	63.6/71.8/86.3	93.8/98.7/93.7	99.7/96.5/91.8
	Cable	70.9/-	78.2/-	57.8/74.0/76.3	84.1/89.5/82.5	95.2/95.9/88.0	55.6/69.7/76.0	55.7/74.8/77.7	94.8/ 98.8/95.2
	Capsule	73.4/-	68.3/-	65.3/92.5/90.4	94.1/96.9/96.9	86.9/ 97.8/94.4	52.9/82.0/90.5	60.5/81.4/90.5	89.0/97.5/95.5
	Hazelnut	85.5/-	97.1/-	93.7/97.5/92.3	60.8/69.8/86.4	99.8/100./99.3	87.0/90.4/88.1	93.0/95.8/89.8	99.5/99.7/97.3
	Metal Nut	88.0/-	64.9/-	72.8/95.0/92.0	100./100./99.5	99.2/99.9/ 99.5	60.0/74.4/89.4	53.0/80.1/89.4	99.1/96.0/91.6
	Pill	68.8/-	79.7/-	82.2/94.9/92.4	97.5/99.6/96.8	93.7/98.7/95.7	55.8/84.0/91.6	62.1/93.1/91.6	95.7/98.5/94.5
	Screw	56.9/-	75.6/-	92.0/95.7/89.9	97.7/99.3/95.8	87.5/96.5/89.0	53.6/71.9/85.9	58.7/81.9/85.6	90.7/ 99.7/97.9
	Toothbrush	95.3/-	75.3/-	90.6/96.8/90.0	97.2/99.0/94.7	94.2/97.4/95.2	57.5/68.0/83.3	78.6/83.9/83.3	99.7/99.9/99.2
	Transistor	86.6/-	73.4/-	74.8/77.4/71.1	94.2/95.2/90.0	99.8/98.0/93.8	57.8/44.6/57.1	61.0/57.8/59.1	99.8/99.6/97.4
	Zipper	79.7/-	87.4/-	98.8/ 99.9/99.2	99.5/99.9/99.2	95.8/99.5/97.1	64.9/77.4/88.1	73.6/89.5/90.6	95.1/99.1/94.4
Textures	Carpet	93.8/-	69.8/-	98.0/99.1/96.7	98.5/99.6/97.2	99.8/99.9/99.4	95.5/98.7/91.0	99.4/99.8/ 99.4	99.4/ 99.9/98.3
	Grid	73.9/-	83.8/-	99.3/99.7/98.2	98.0/99.4/96.5	98.2/99.5/97.3	83.5/93.9/86.9	67.3/82.6/84.4	98.5/ 99.8/97.7
	Leather	99.9/-	93.6/-	98.7/99.3/95.0	100./100./100.	100./100./100.	98.4/99.5/96.3	97.4/99.0/96.3	99.8/99.7/97.6
	Tile	93.3/-	89.5/-	99.8/100./100.	98.3/99.3/96.4	99.3/99.8/98.2	93.6/97.5/92.0	97.1/98.7/94.1	96.8/99.9/98.4
	Wood	98.4/-	93.4/-	99.8/100./100.	99.2/99.8/98.3	98.6/99.6/96.6	98.6/99.6/97.5	97.8/99.4/95.9	99.7/ 100./100.
Mean	84.2/-	81.9/-	88.1/94.7/92.0	94.6/96.5/95.2	96.5/98.8/96.2	71.9/81.6/86.6	76.6/87.8/88.1	97.2/99.0/96.5	

Table 1: Comparison with SOTA methods on MVTec-AD dataset for multi-class anomaly detection with $AUROC_{cls}/AP_{cls}/F1max_{cls}$ metrics.

As Batch Normalization (BN) (Ioffe and Szegedy 2015) considers the normalization statistics of all images within a batch, it leads to a loss of unique details in each sample. BN is suitable for a relatively large mini-batch scenario with similar data distributions. However, for multi-class anomaly detection where there are significant differences in data distributions among different categories, normalizing the entire batch is not suitable for tasks in the multi-class setting. Since the results generated by using SD mainly depend on the input image instance, using Instance Normalization (IN) (Ulyanov, Vedaldi, and Lempitsky 2017) can not only accelerate model convergence but also maintain the independence between each image instance. In addition, in terms of choosing the activation function, we use the SiLU (Elfwing, Uchibe, and Doya 2018) instead of the commonly used ReLU (Hahnloser et al. 2000), which can preserve more input information. Experimental results in Table 7 show that the performance is improved by using IN and SiLU simultaneously instead of the combination of BN and ReLU.

Anomaly localization and detection

During the inference stage, the reconstruction image is obtained through the diffusion and denoising process in the latent space. For anomaly localization and detection, We use the same ImageNet pre-trained feature extractor Ψ to extract features from both the input image x_0 and the reconstructed image \hat{x}_0 and calculate the anomaly map on different scale feature maps \mathcal{M}^n using cosine similarity:

$$\mathcal{M}^n(x_0, \hat{x}_0) = 1 - \frac{(\Psi^n(x_0, \hat{x}_0))^T \cdot \Psi^n(x_0, \hat{x}_0)}{\|\Psi^n(x_0, \hat{x}_0)\| \|\Psi^n(x_0, \hat{x}_0)\|}, \quad (9)$$

where n represents the n -th feature layer f_n and the anomaly score \mathcal{S} for an input-pair of anomaly localization is:

$$\mathcal{S} = \sum_{n \in N} \sigma_n \mathcal{M}^n(x_0, \hat{x}_0), \quad (10)$$

Metrics	Non-Diffusion		Diffusion-based		
	DRAEM	UniAD	DDPM	LDM	Ours
$AUROC_{cls}$	79.1	85.5	54.5	56.7	86.8
AP_{cls}	81.9	85.5	57.9	61.4	88.3
$F1max_{cls}$	78.9	84.4	72.3	73.1	85.1
$AUROC_{seg}$	91.3	95.9	79.7	86.6	96.0
AP_{seg}	23.5	21.0	2.2	6.0	26.1
$F1max_{seg}$	29.5	27.0	4.5	9.9	33.0
PRO	58.8	75.6	46.8	55.0	75.2

Table 2: Quantitative comparisons on VisA dataset.

where σ_n indicates the upsampling factor in order to keep the same dimension of the pixel space image and N indicates the number of feature layers used during inference.

Experiment

Datasets and evaluation metrics

MVTec-AD dataset. MVTEC-AD (Bergmann et al. 2019) dataset simulates real-world industrial production scenarios, filling the gap in unsupervised anomaly detection. It consists of 5 types of textures and 10 types of objects, in 5,354 high-resolution images from different domains. The training set contains 3,629 images with only anomaly-free samples. The test set consists of 1,725 images, including both normal and abnormal samples. Pixel-level annotations are provided for the anomaly localization evaluation.

VisA dataset. VisA (Zou et al. 2022) dataset consists of a total of 10,821 high-resolution images, including 9,621 normal images and 1,200 anomaly images with 78 types of anomalies. The VisA dataset comprises 12 subsets, each corresponding to a distinct object. 12 objects could be categorized into three different object types: Complex structure, Multiple instances, and Single instance.

MVTec-3D dataset. MVTEC-3D (Bergmann et al. 2022) dataset comprises 4,147 scans obtained using a high-resolution industrial 3D sensor. It consists of 10 categories

Category		Non-Diffusion Method					Diffusion-based Method		
		PaDiM	MKD	DRAEM	RD4AD	UniAD	DDPM	LDM	Ours
Objects	Bottle	96.1/-	91.8/-	87.6/62.5/56.9	97.8/ 68.2 /67.6	98.1/66.0/ 69.2	59.9/ 4.9/11.7	86.9/49.1/50.0	98.4 /52.2/54.8
	Cable	81.0/-	89.3/-	71.3/14.7/17.8	85.1/26.3/33.6	97.3/39.9/45.2	66.5/ 6.7/10.6	89.3/18.5/26.2	96.8 / 50.1 / 57.8
	Capsule	96.9/-	88.3/-	50.5/ 6.0/10.0	98.8 / 43.4 / 50.0	98.5/42.7/46.5	63.1/ 6.2/ 9.7	90.0/ 7.9/27.3	97.1/42.0/45.3
	Hazelnut	96.3/-	91.2/-	96.9/70.0/60.5	97.9/36.2/51.6	98.1/55.2/56.8	91.2/24.1/28.3	95.1/51.2/53.5	98.3 / 79.2 / 80.4
	Metal Nut	84.8/-	64.2/-	62.2/31.1/21.0	93.8/ 62.3 /65.4	94.8/55.5/ 66.4	62.7/14.6/29.2	70.5/19.3/30.7	97.3 /30.0/38.3
	Pill	87.7/-	69.7/-	94.4/59.1/44.1	97.5 / 63.4 / 65.2	95.0/44.0/53.9	55.3/ 4.0/ 8.4	74.9/10.2/15.0	95.7/46.0/51.4
	Screw	94.1/-	92.1/-	95.5/33.8/40.6	99.4 /40.2/44.6	98.3/28.7/37.6	91.1/ 1.8/ 3.8	91.7/ 2.2/ 4.6	97.9/ 60.6 / 59.6
	Toothbrush	95.6/-	88.9/-	97.7/55.2/55.8	99.0 /53.6/58.8	98.4/34.9/45.7	76.9/ 4.0/ 7.7	93.7/20.4/ 9.8	99.0 / 78.7 / 72.8
	Transistor	92.3/-	71.7/-	64.5/23.6/15.1	85.9/42.3/45.2	97.9 / 59.5 / 64.6	53.2/ 5.8/11.4	85.5/25.0/30.7	95.1/15.6/31.7
	Zipper	94.8/-	86.1/-	98.3/ 74.3 / 69.3	98.5 /53.9/60.3	96.8/40.1/49.9	67.4/ 3.5/ 7.6	66.9/ 5.3/ 7.4	96.2/60.7/60.0
Textures	Carpet	97.6/-	95.5/-	98.6/ 78.7 / 73.1	99.0/58.5/60.4	98.5/49.9/51.1	89.2/18.8/44.3	99.1 /70.6/66.0	98.6/42.2/46.4
	Grid	71.0/-	82.3/-	98.7/44.5/46.2	99.2 /46.0/47.4	96.5/23.0/28.4	63.1/ 0.7/ 1.9	52.4/ 1.1/ 1.9	96.6 / 66.0 / 64.1
	Leather	84.8/-	96.7/-	97.3/ 60.3 /57.4	99.3 /38.0/45.1	98.8/32.9/34.4	97.3/38.9/43.2	99.0/45.9/44.0	98.8/56.1/ 62.3
	Tile	80.5/-	85.3/-	98.0 / 93.6 / 86.0	95.3/48.5/60.5	91.8/42.1/50.6	87.0/35.2/36.6	90.1/43.9/51.6	92.4/65.7/64.1
	Wood	89.1/-	80.5/-	96.0 / 81.4 / 74.6	95.3/47.8/51.0	93.2/37.2/41.5	84.7/30.9/37.3	92.3/44.1/46.6	93.3/43.3/43.5
Mean	89.5/-	84.9/-	87.2/52.5/48.6	96.1/48.6/53.8	96.8 /43.4/49.5	75.6/13.3/19.5	85.1/27.6/31.0	96.8 / 52.6 / 55.5	

Table 3: Comparison with SOTA methods on MVTec-AD dataset for multi-class anomaly localization with $AUROC_{seg}/AP_{seg}/F1max_{seg}$ metrics.

Method	Non-Diffusion		Diffusion-based		
	DRAEM	UniAD	DDPM	LDM	Ours
PRO	71.1	90.4	49.0	66.3	90.7

Table 4: Multi-class anomaly localization results with PRO metric on MVTec-AD datasets.

with both RGB images and 3D point clouds respectively. The training set contains 2,656 images with only anomaly-free samples. The test set consists of 1,197 images, including both normal and abnormal samples. Only RGB images are used in this experiment.

Medical dataset. We also merge three types of medical datasets BraTS2021 (Baid et al. 2021), BTCV (Landman et al. 2015) and LiTs (Bilic et al. 2023) into one *Medical* dataset for multi-class anomaly detection. The training set contains 9,042 slices and the test set consists of 5,208 slices.

Evaluation Metrics. Following prior works, Area Under the Receiver Operating Characteristic Curve (AUROC), Average Precision (AP) and F1-score-max (F1max) are used in both anomaly detection and anomaly localization, where *cls* represents the image level anomaly detection and *seg* represents the pixel level anomaly localization. Also, Per-Region-Overlap (PRO) is used in anomaly localization. The DICE score is commonly used in the medical field.

Implementation Details

All images in MVTec-AD and VisA are resized to 256×256 . For the denoising network, we adopt the 4-th block of SGDB for connection to SDDb. In this experiment, we adopt ResNet50 as the feature extraction network and choose $n \in \{2, 3, 4\}$ as the feature layers used in calculating the anomaly localization. We utilized the KL method as the Auto-encoder and fine-tune the model before training the denoising network. We train for 1000 epochs on a single NVIDIA Tesla V100 32GB with a batch size of 12. Adam

optimiser (Loshchilov and Hutter 2019) with a learning rate of $1e^{-5}$ is set. A Gaussian filter with $\sigma = 5$ is used to smooth the anomaly localization score. For anomaly detection, the anomaly score of the image is the maximum value of the averagely pooled anomaly localization score which undergoes 8 rounds of global average pooling operations with a size of 8×8 . During inference, the initial denoising timestep T is set from 1,000. We use DDIM (Song, Meng, and Ermon 2021) as the sampler with 10 steps by default.

Comparison with SOTAs

We conduct and analyze a range of qualitative and quantitative comparison experiments on MVTec-AD, VisA, MVTec-3D and *Medical* datasets. We choose a synthesizing-based method DRAEM (Zavrtanik, Kristan, and Skočaj 2021a), three embedding-based methods MKD (Salehi et al. 2021), PaDiM (Defard et al. 2021) and RD4AD (Deng and Li 2022), a reconstruction-based method EdgRec (Liu et al. 2022), a unified SOTA UniAD (You et al. 2022) method and diffusion-based DDPM and LDM methods. Specifically, we categorize the aforementioned methods into two types: non-diffusion and diffusion-based methods. For the experiments on *Medical* dataset, we follow the BMAD (Bao et al. 2023) benchmark and add two methods STFPM (Yamada and Hotta 2021) and CFLOW (Gudovskiy, Ishizaka, and Kozuka 2022) for comparison.

Qualitative Results. We conducted substantial qualitative experiments on MVTec-AD and VisA datasets to visually demonstrate the superiority of our method in image reconstruction and the accuracy of anomaly localization. As shown in Figure 4, our method exhibits better reconstruction capabilities for anomalous regions compared to the EdgRec on MVTec-AD dataset. In comparison to UniAD shown in Figure 5, our method exhibits more accurate anomaly localization abilities on VisA dataset. More qualitative results will be presented in *Appendix*.

Quantitative Results. As shown in Table 1 and in Ta-

Metrics	Non-Diffusion		Diffusion-based		
	DRAEM	UniAD	DDPM	LDM	Ours
$AUROC_{cls}$	63.2	78.9	66.3	68.5	84.6
AP_{cls}	86.1	93.4	78.0	90.6	94.8
$F1max_{cls}$	89.2	91.4	86.6	91.6	95.5
$AUROC_{seg}$	93.2	96.5	90.7	92.2	96.4
AP_{seg}	16.8	21.2	6.0	9.3	25.3
$F1max_{seg}$	20.2	28.0	10.7	13.5	32.2
PRO	55.0	88.1	69.7	73.8	87.8

Table 5: Quantitative comparisons on MVTec-3D dataset.

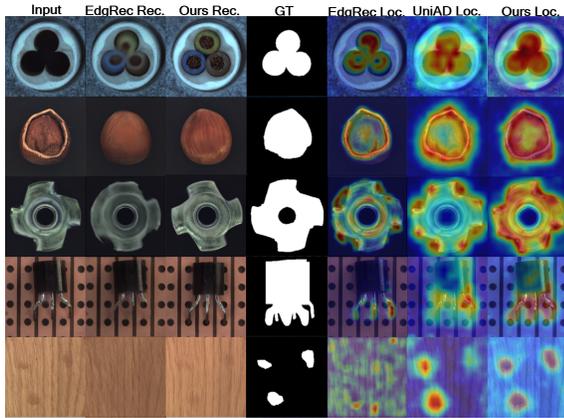


Figure 4: Qualitative illustration on MVTec-AD dataset.

ble 3, our method achieves SOTA AUROC/AP/F1max metrics of 97.2/99.0/96.5 and 96.8/52.6/55.5 for image-wise and pixel-wise respectively for multi-class setting on MVTec-AD dataset. For the diffusion-based methods, our approach significantly outperforms existing DDPM and LDM methods in terms of 11.7 \uparrow in AUROC and 25 \uparrow in AP for anomaly localization. For non-diffusion methods, our approach surpasses existing methods in both metrics, especially at the pixel level, where our method exceeds UniAD by 9.2 \uparrow /6.0 \uparrow in AP/F1max. Our method has also demonstrated its superiority on VisA dataset, as shown in Table 2. Our approach exhibits significant improvements compared to diffusion-based methods of 30.1 \uparrow /9.4 \uparrow than the LDM method in image/pixel AUROC. It also performs well compared to UniAD by 4.9 \uparrow /6.0 \uparrow in pixel AP/F1max metrics. Detailed experiments for each category are provided in *Appendix*. We have extended the method to 3D datasets and medical domain datasets. Table 5 and Table 6 show the effectiveness and scalability of our method on MVtec-3D and *Medical* datasets, with results surpassing the state of the art (SOTA).

Metrics	MKD	CFLOW	RD4AD	PaDiM	PatchCore	STFPM	UniAD	Ours
$AUROC_{cls}$	70.9	62.0	74.7	64.6	76.0	72.2	76.4	77.2
$AUROC_{seg}$	92.8	93.2	96.2	93.0	96.8	93.4	96.7	96.9
PRO	79.3	79.0	88.0	79.2	86.6	86.0	87.4	87.7
$DICE$	21.9	13.5	19.5	15.2	21.7	17.1	28.7	32.3

Table 6: Quantitative comparisons on *Medical* dataset.

SD	MSG	SGBE3	SGBE4	BN+ReLU	IN+SiLU	cls	seg
✓						79.3	89.5
✓	✓					95.1	91.1
✓	✓	✓				95.3	89.1
✓	✓	✓	✓			93.8	91.2
✓	✓	✓		✓		96.7	96.7
✓	✓	✓			✓	97.2	96.8

Table 7: Ablation studies on the design of DiAD with AUROC metrics.

Ablation Studies

The architecture design of DiAD. We investigate the importance of each module in DiAD as shown in Table 7. SD indicates only the diffusion model without connecting to the SG network which is the LDM’s architecture. MSG indicates only the middle block of the SG network adding to the middle of SD. SGBE3 and SGBE4 indicate directly skip-connecting to the corresponding SDDB. When connecting SGBE3 and SGBE4 at the same time, more details of the original images are preserved in terms of texture, but the reconstruction ability for large anomaly areas decreases. Using the combination of IN+SiLU in the SFF block yields better results compared to using BN+ReLU.

Effect of pre-trained feature extractors. Table 8 shows the quantitative comparison of using different pre-trained backbones as feature extraction networks. ResNet50 achieved the best performance in anomaly classification metrics, while WideResNet101 excelled in anomaly segmentation.

Backbone		$AUROC_{cls}$	AP_{cls}	$F1max_{cls}$	$AUROC_{seg}$	AP_{seg}	$F1max_{seg}$	PRO
VGG	16	91.8	97.2	93.9	92.1	47.2	50.5	80.1
	19	91.3	96.9	93.7	92.3	47.5	50.6	80.4
ResNet	18	94.7	98.1	96.0	96.0	49.9	53.3	89.1
	34	95.2	98.3	95.7	96.2	51.2	54.5	89.6
	50	97.2	99	96.5	96.8	52.6	55.5	90.7
	101	96.2	98.4	96.5	96.9	52.9	56.4	91.2
WideResNet	50	95.9	98.6	96.5	96.4	51.8	55.1	89.3
	101	95.6	98.3	95.8	96.9	54.6	56.5	91.4
EfficientNet	b0	93.5	97.7	94.7	94.0	50.0	52.4	84.0
	b2	94.2	98.0	95.1	94.1	48.6	52.1	84.2
	b4	92.8	97.5	94.8	93.6	47.2	50.7	83.5

Table 8: Ablation studies on different feature extractors.

Effect of feature layers used in anomaly score calculation. After extracting feature maps of 5 different scales using a pre-trained backbone, the anomaly scores are calculated by computing the cosine similarity between feature maps from different layers. The experimental results, as shown in Table 9, indicate that using feature maps from layers f_2 , f_3 , and f_4 (with corresponding sizes of 64×64 , 32×32 , and 16×16) yields the best performance.

f_1	f_2	f_3	f_4	f_5	$AUROC_{cls}$	AP_{cls}	$F1max_{cls}$	$AUROC_{seg}$	AP_{seg}	$F1max_{seg}$
✓	✓	✓	✓	✓	93.8	97.8	95.0	94.0	42.0	45.9
✓	✓	✓	✓		96.7	98.7	96.1	96.7	52.5	55.2
✓	✓	✓	✓		93.4	97.1	93.6	95.2	48.5	51.3
✓	✓	✓	✓		97.1	99.0	96.8	96.4	49.4	53.1
✓	✓	✓	✓		97.2	99.0	96.5	96.8	52.6	55.5
✓	✓	✓			94	97.4	94.2	95.3	48.5	51.7
✓	✓	✓			97.1	99.0	96.8	96.4	49.4	53.1

Table 9: Ablation studies on the feature layers used in calculating the anomaly localization score based on ResNet50.

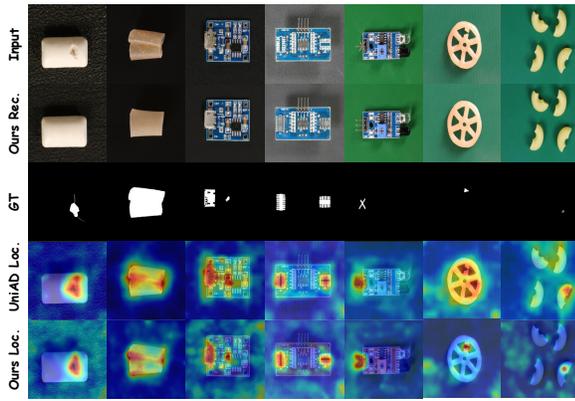


Figure 5: Qualitative results on VisA dataset.

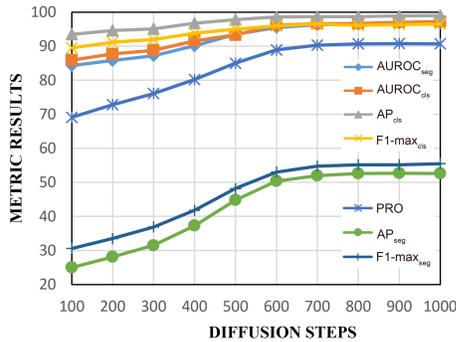


Figure 6: Ablation studies on different diffusion timesteps.

Effect of forward diffusion timesteps. Increasing the number of diffusion steps in the forward process impacts the performance of image reconstruction. The experimental results, depicted in Figure 6, indicate that with an increasing number of forward diffusion steps, the image approaches pure Gaussian noise, while the anomaly reconstruction ability improves as well. Nevertheless, when the number of forward diffusion steps is less than 600, a significant decline in performance occurs because the number of steps is insufficient for anomaly reconstruction.

Conclusion

This paper proposes a diffusion-based DiAD framework to address the issue of category and semantic loss in the stable diffusion model for multi-class anomaly detection. We propose the Semantic-Guided network and Spatial-aware Feature Fusion block to better reconstruct the abnormal regions while maintaining the same semantic information as the input image. Our approach achieves state-of-the-art performance on MVTEC-AD and VisA datasets, significantly outperforming the non-diffusion and diffusion-based methods.

Limitation. Although our method has demonstrated exceptional performance in reconstructing anomalies, it can be susceptible to the influence of background impurities, resulting in errors in localization and classification. In the future, we will further explore diffusion models and enhance

the background’s anti-interference capability for multi-class anomaly detection. Additionally, we will incorporate multi-modal assistance in our anomaly detection. Lastly, we will utilize larger models to enhance reconstruction performance.

References

- Amit, T.; Shaharbany, T.; Nachmani, E.; and Wolf, L. 2022. SegDiff: Image Segmentation with Diffusion Probabilistic Models. *arXiv:2112.00390*.
- Baid, U.; Ghodasara, S.; Mohan, S.; Bilello, M.; Calabrese, E.; Colak, E.; Farahani, K.; Kalpathy-Cramer, J.; Kitamura, F. C.; Pati, S.; et al. 2021. The rsna-asnr-miccai brats 2021 benchmark on brain tumor segmentation and radiogenomic classification. *arXiv preprint arXiv:2107.02314*.
- Bao, J.; Sun, H.; Deng, H.; He, Y.; Zhang, Z.; and Li, X. 2023. BMAD: Benchmarks for Medical Anomaly Detection. *arXiv preprint arXiv:2306.11876*.
- Bergmann, P.; Fauser, M.; Sattlegger, D.; and Steger, C. 2019. MVTEC AD—A comprehensive real-world dataset for unsupervised anomaly detection. In *CVPR*, 9592–9600.
- Bergmann, P.; Jin, X.; Sattlegger, D.; and Steger, C. 2022. The MVTEC 3D-AD Dataset for Unsupervised 3D Anomaly Detection and Localization. In *VISGRAPP*. SCITEPRESS - Science and Technology Publications.
- Bilic, P.; Christ, P.; Li, H. B.; Vorontsov, E.; Ben-Cohen, A.; Kaissis, G.; Szeskin, A.; Jacobs, C.; Mamani, G. E. H.; Chartrand, G.; et al. 2023. The liver tumor segmentation benchmark (lits). *Medical Image Analysis*, 84: 102680.
- Cao, Y.; Wan, Q.; Shen, W.; and Gao, L. 2022. Informative knowledge distillation for image anomaly segmentation. *Knowledge-Based Systems*, 248: 108846.
- Cao, Y.; Xu, X.; Sun, C.; Cheng, Y.; Du, Z.; Gao, L.; and Shen, W. 2023. Segment Any Anomaly without Training via Hybrid Prompt Regularization. *arXiv preprint arXiv:2305.10724*.
- Chen, R.; Xie, G.; Liu, J.; Wang, J.; Luo, Z.; Wang, J.; and Zheng, F. 2023a. Easynet: An easy network for 3d industrial anomaly detection. In *ACM MM*, 7038–7046.
- Chen, S.; Sun, P.; Song, Y.; and Luo, P. 2022. DiffusionDet: Diffusion Model for Object Detection. *arXiv:2211.09788*.
- Chen, X.; Han, Y.; and Zhang, J. 2023. A Zero-/Few-Shot Anomaly Classification and Segmentation Method for CVPR 2023 VAND Workshop Challenge Tracks 1&2: 1st Place on Zero-shot AD and 4th Place on Few-shot AD. *arXiv preprint arXiv:2305.17382*.
- Chen, X.; Zhang, J.; Tian, G.; He, H.; Zhang, W.; Wang, Y.; Wang, C.; Wu, Y.; and Liu, Y. 2023b. CLIP-AD: A Language-Guided Staged Dual-Path Model for Zero-shot Anomaly Detection. *arXiv preprint arXiv:2311.00453*.
- Defard, T.; Setkov, A.; Loesch, A.; and Audigier, R. 2021. Padim: a patch distribution modeling framework for anomaly detection and localization. In *ICPR*, 475–489. Springer.
- Deng, H.; and Li, X. 2022. Anomaly detection via reverse distillation from one-class embedding. In *CVPR*, 9737–9746.

- Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale hierarchical image database. In *CVPR*, 248–255. Ieee.
- Ding, C.; Pang, G.; and Shen, C. 2022. Catching both gray and black swans: Open-set supervised anomaly detection. In *CVPR*, 7388–7398.
- Elfwing, S.; Uchibe, E.; and Doya, K. 2018. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. *Neural networks*, 107: 3–11.
- Gu, Z.; Liu, L.; Chen, X.; Yi, R.; Zhang, J.; Wang, Y.; Wang, C.; Shu, A.; Jiang, G.; and Ma, L. 2023. Remembering Normality: Memory-guided Knowledge Distillation for Unsupervised Anomaly Detection. In *ICCV*, 16401–16409.
- Gudovskiy, D.; Ishizaka, S.; and Kozuka, K. 2022. Cflowad: Real-time unsupervised anomaly detection with localization via conditional normalizing flows. In *WACV*, 98–107.
- Hahnloser, R. H.; Sarpeshkar, R.; Mahowald, M. A.; Douglas, R. J.; and Seung, H. S. 2000. Digital selection and analogue amplification coexist in a cortex-inspired silicon circuit. *nature*, 405(6789): 947–951.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *CVPR*, 770–778.
- Ho, J.; Chan, W.; Saharia, C.; Whang, J.; Gao, R.; Gritsenko, A.; Kingma, D. P.; Poole, B.; Norouzi, M.; Fleet, D. J.; and Salimans, T. 2022. Imagen Video: High Definition Video Generation with Diffusion Models. arXiv:2210.02303.
- Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising Diffusion Probabilistic Models. In *NeurIPS*, volume 33, 6840–6851.
- Huang, C.; Guan, H.; Jiang, A.; Zhang, Y.; Spratling, M.; and Wang, Y.-F. 2022. Registration based few-shot anomaly detection. In *ECCV*, 303–319. Springer.
- Ioffe, S.; and Szegedy, C. 2015. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. In Bach, F. R.; and Blei, D. M., eds., *ICML*, volume 37 of *JMLR Workshop and Conference Proceedings*, 448–456. JMLR.org.
- Jeong, J.; Zou, Y.; Kim, T.; Zhang, D.; Ravichandran, A.; and Dabeer, O. 2023. Winclip: Zero-/few-shot anomaly classification and segmentation. In *CVPR*, 19606–19616.
- Kingma, D. P.; and Welling, M. 2022. Auto-Encoding Variational Bayes. arXiv:1312.6114.
- Landman, B.; Xu, Z.; Igelsias, J.; Styner, M.; Langerak, T.; and Klein, A. 2015. Miccai multi-atlas labeling beyond the cranial vault—workshop and challenge. In *Proc. MICCAI Multi-Atlas Labeling Beyond Cranial Vault—Workshop Challenge*, volume 5, 12.
- Li, C.-L.; Sohn, K.; Yoon, J.; and Pfister, T. 2021. Cutpaste: Self-supervised learning for anomaly detection and localization. In *CVPR*, 9664–9674.
- Liang, Y.; Zhang, J.; Zhao, S.; Wu, R.; Liu, Y.; and Pan, S. 2023. Omni-frequency channel-selection representations for unsupervised anomaly detection. *IEEE Transactions on Image Processing*.
- Liu, J.; Xie, G.; Wang, J.; Li, S.; Wang, C.; Zheng, F.; and Jin, Y. 2023. Deep Industrial Image Anomaly Detection: A Survey. *arXiv preprint arXiv:2301.11514*, 2.
- Liu, T.; Li, B.; Zhao, Z.; Du, X.; Jiang, B.; and Geng, L. 2022. Reconstruction from edge image combined with color and gradient difference for industrial surface anomaly detection. arXiv:2210.14485.
- Liznerski, P.; Ruff, L.; Vandermeulen, R. A.; Franks, B. J.; Kloft, M.; and Müller, K. 2021. Explainable Deep One-Class Classification. In *ICLR*.
- Loshchilov, I.; and Hutter, F. 2019. Decoupled Weight Decay Regularization. arXiv:1711.05101.
- Mousakhan, A.; Brox, T.; and Tayyub, J. 2023. Anomaly Detection with Conditioned Denoising Diffusion Models. arXiv:2305.15956.
- Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-Resolution Image Synthesis with Latent Diffusion Models. arXiv:2112.10752.
- Ronneberger, O.; Fischer, P.; and Brox, T. 2015. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 234–241. Springer.
- Roth, K.; Pemula, L.; Zepeda, J.; Schölkopf, B.; Brox, T.; and Gehler, P. 2022. Towards total recall in industrial anomaly detection. In *CVPR*, 14318–14328.
- Salehi, M.; Mirzaei, H.; Hendrycks, D.; Li, Y.; Rohban, M. H.; and Sabokrou, M. 2022. A Unified Survey on Anomaly, Novelty, Open-Set, and Out-of-Distribution Detection: Solutions and Future Challenges. arXiv:2110.14051.
- Salehi, M.; Sadjadi, N.; Baselizadeh, S.; Rohban, M. H.; and Rabiee, H. R. 2021. Multiresolution knowledge distillation for anomaly detection. In *CVPR*, 14902–14912.
- Song, J.; Meng, C.; and Ermon, S. 2021. Denoising Diffusion Implicit Models. In *ICLR*. OpenReview.net.
- Tan, D. S.; Chen, Y.-C.; Chen, T. P.-C.; and Chen, W.-C. 2021. Trustmae: A noise-resilient defect classification framework using memory-augmented auto-encoders with trust regions. In *WACV*, 276–285.
- Tan, M.; and Le, Q. 2019. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, 6105–6114. PMLR.
- Tao, X.; Gong, X.; Zhang, X.; Yan, S.; and Adak, C. 2022. Deep Learning for Unsupervised Anomaly Localization in Industrial Images: A Survey. *IEEE Transactions on Instrumentation and Measurement*, 71: 1–21.
- Ulyanov, D.; Vedaldi, A.; and Lempitsky, V. 2017. Instance Normalization: The Missing Ingredient for Fast Stylization. arXiv:1607.08022.
- Wang, Y.; Peng, J.; Zhang, J.; Yi, R.; Wang, Y.; and Wang, C. 2023. Multimodal Industrial Anomaly Detection via Hybrid Fusion. In *CVPR*, 8032–8041.
- Wu, J.; Li, J.; Zhang, J.; Zhang, B.; Chi, M.; Wang, Y.; and Wang, C. 2023. PVG: Progressive Vision Graph for Vision Recognition. *arXiv preprint arXiv:2308.00574*.

Wyatt, J.; Leach, A.; Schmon, S. M.; and Willcocks, C. G. 2022. AnoDDPM: Anomaly Detection with Denoising Diffusion Probabilistic Models using Simplex Noise. In *CVPR Workshops 2022, New Orleans, LA, USA, June 19-20, 2022*, 649–655. IEEE.

Xie, G.; Wang, J.; Liu, J.; Jin, Y.; and Zheng, F. 2023. Pushing the Limits of Fewshot Anomaly Detection in Industry Vision: Graphcore. In *ICLR*.

Yamada, S.; and Hotta, K. 2021. Reconstruction student with attention for student-teacher pyramid matching. *arXiv preprint arXiv:2111.15376*.

Yan, X.; Zhang, H.; Xu, X.; Hu, X.; and Heng, P. 2021. Learning Semantic Context from Normal Samples for Unsupervised Anomaly Detection. In *AAAI*, 3110–3118.

Yi, J.; and Yoon, S. 2020. Patch SVDD: Patch-level SVDD for Anomaly Detection and Segmentation. In *ACCV*.

Yoon, J.; Sohn, K.; Li, C.-L.; Arik, S. O.; Lee, C.-Y.; and Pfister, T. 2022. Self-supervise, Refine, Repeat: Improving Unsupervised Anomaly Detection. *Transactions on Machine Learning Research*.

You, Z.; Cui, L.; Shen, Y.; Yang, K.; Lu, X.; Zheng, Y.; and Le, X. 2022. A Unified Model for Multi-class Anomaly Detection. In *NeurIPS*, volume 35, 4571–4584.

Yu, J.; Zheng, Y.; Wang, X.; Li, W.; Wu, Y.; Zhao, R.; and Wu, L. 2021. FastFlow: Unsupervised Anomaly Detection and Localization via 2D Normalizing Flows. *arXiv:2111.07677*.

Zagoruyko, S.; and Komodakis, N. 2016. Wide Residual Networks. In *BMVC*. BMVA Press.

Zavrtnik, V.; Kristan, M.; and Skočaj, D. 2021a. Draem-a discriminatively trained reconstruction embedding for surface anomaly detection. In *ICCV*, 8330–8339.

Zavrtnik, V.; Kristan, M.; and Skočaj, D. 2021b. Reconstruction by inpainting for visual anomaly detection. *Pattern Recognition*, 112: 107706.

Zhang, H.; Wang, Z.; Wu, Z.; and Jiang, Y.-G. 2023a. DiffusionAD: Denoising Diffusion for Anomaly Detection. *arXiv:2303.08730*.

Zhang, J.; Chen, X.; Xue, Z.; Wang, Y.; Wang, C.; and Liu, Y. 2023b. Exploring Grounding Potential of VQA-oriented GPT-4V for Zero-shot Anomaly Detection. *arXiv preprint arXiv:2311.02612*.

Zhang, J.; Li, X.; Li, J.; Liu, L.; Xue, Z.; Zhang, B.; Jiang, Z.; Huang, T.; Wang, Y.; and Wang, C. 2023c. Rethinking Mobile Block for Efficient Attention-based Models. In *ICCV*, 1389–1400.

Zhang, J.; Li, X.; Wang, Y.; Wang, C.; Yang, Y.; Liu, Y.; and Tao, D. 2022. Eatformer: Improving vision transformer inspired by evolutionary algorithm. *arXiv preprint arXiv:2206.09325*.

Zhang, L.; and Agrawala, M. 2023. Adding Conditional Control to Text-to-Image Diffusion Models. *arXiv:2302.05543*.

Zou, Y.; Jeong, J.; Pemula, L.; Zhang, D.; and Dabeer, O. 2022. Spot-the-difference self-supervised pre-training for

anomaly detection and segmentation. In *ECCV*, 392–408. Springer.

Appendices

Effect of DDIM sampler steps

In order to accelerate the sampling speed in the denoising process, UiAD adopts the DDIM sampling strategy. We investigated the impact of different DDIM sampler steps on the results, as shown in Table 10. The results indicate that increasing the number of sampling steps does not significantly affect the results. Therefore, using a 10-step sampling process can achieve the best performance while greatly accelerating the sampling speed.

Steps	1	5	10	20	50	100	200
<i>seg</i>	72.5	96.5	96.8	96.8	96.7	96.7	96.8
<i>cls</i>	66.1	96.4	97.2	97.1	97.0	96.8	96.9

Table 10: Ablation studies on DDIM sampler steps.

Effect of Global average pooling

Global average pooling is used to reduce the potential occurrence of false positives. For m - n in the table below, m represents the iterations and n represents the kernel size. Through quantitative analysis, the most effective approach is employing an 8×8 size global average pooling with 8 iterations. Also, the best-performing combinations exhibit the same feature map size.

Global Average Pooling	1-16	4-16	5-12	6-10	8-8	10-7	15-5	20-4
AUROC-cls	96.0	96.7	96.9	97.1	97.2	97.2	97.0	96.8

Limitations of the datasets

We found that there are several categories of image-level anomaly detection results that are significantly lower than others, such as capsules and screws. As shown in Fig 7, we discovered some false positives in input good images during the test. Our method performs well in reconstructing the objects in the objects’ main bodies, but the background region of the original image contains impurities, causing the pre-trained feature extraction network to extract features that perceive the background impurities as anomalies. As anomaly detection is expected to identify anomalies within the object rather than the background region, there are certain deficiencies in the Mvtec-AD as well as the VisA datasets that lead to false positives. In response to this issue, we increase the number of global average pooling operations to alleviate the problem of high anomaly scores caused by impurities in the background.

Hyperparameters of DiAD

We provided a comprehensive set of hyperparameters for the three models in DiAD as shown in Table 11.

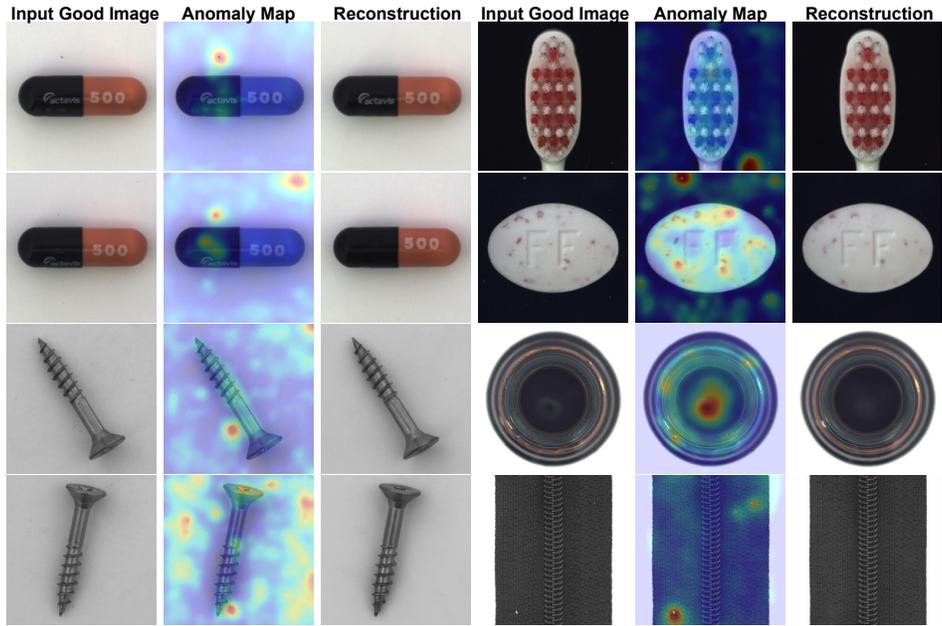


Figure 7: Visualization of false positive classifications and localizations.

Parameters Name	Model Name		
	SD Denoising Network	SG Network	Autoencoder
z shape	$32 \times 32 \times 4$		
$ z $	4096		
Diffusion steps T	1000		
DDIM sampling steps T	10		
Noise Schedule	linear		
Model input shape	$32 \times 32 \times 4$	$256 \times 256 \times 3$	$256 \times 256 \times 3$
N params	859M	471M	83.7M
Embed dim	-	-	4
Channels	320	320	128
Num res blocks	2	2	2
Channel Multiplier	1,2,4,4	1,2,4,4	1,2,4,4
Attention resolutions	4,2,1	4,2,1	-
Num Heads	8	8	-
Batch Size	12		
Accumulate_grad_batches	4		
Epochs	1000		
Learning Rate	1.0e-5		

Table 11: Hyperparameters for the DiAD. All models trained on a single NVIDIA Tesla V100 32GB.

Category	Non-Diffusion Method		Diffusion-based Method		
	DRAEM	UniAD	DDPM	LDM	Ours
pcb1	71.9/72.2/70.0	92.8/92.7/87.8	54.1/47.7/67.1	51.2/46.9/66.8	88.1/88.7/80.7
pcb2	78.4/78.2/76.2	87.8/87.7/83.1	50.8/48.5/66.6	57.0/63.4/67.5	91.4/91.4/84.7
pcb3	76.6/77.4/74.7	78.6/78.6/76.1	53.4/51.2/66.8	62.7/69.6/72.0	86.2/87.6/77.6
pcb4	97.3/97.5/93.5	98.8/98.8/94.3	56.0/48.4/66.4	54.4/47.1/66.8	99.6/99.5/97.0
macaroni1	69.8/68.5/70.9	79.9/79.8/72.7	50.9/55.1/68.0	56.2/49.6/68.4	85.7/85.2/78.8
macaroni2	59.4/60.7/68.0	71.6/71.6/69.9	54.4/51.8/67.1	56.8/52.7/66.6	62.5/57.4/69.6
capsules	83.4/91.1/82.1	55.6/55.6/76.9	58.9/62.7/78.2	57.7/71.4/77.3	58.2/69.0/78.5
candle	69.3/73.9/68.0	94.1/94.0/86.1	52.7/48.3/66.6	50.4/52.2/68.2	92.8/92.0/87.6
cashew	81.7/89.7/87.3	92.8/92.8/91.4	63.5/78.9/80.6	61.1/71.0/80.0	91.5/95.7/89.7
chewinggum	93.7/97.1/91.0	96.3/96.2/95.2	50.9/65.6/80.0	53.9/65.8/81.3	99.1/99.5/95.9
fryum	89.1/95.0/86.6	83.0/83.0/85.0	51.0/62.4/80.0	63.5/71.6/81.6	89.8/95.0/87.2
pipe.fryum	82.8/91.2/83.9	94.7/94.7/93.9	56.9/74.9/80.0	56.1/75.5/80.3	96.2/98.1/93.7
Mean	79.1/81.9/78.9	85.5/85.5/84.4	54.5/57.9/72.3	56.7/61.4/73.1	86.8/88.3/85.1

Table 12: Comparison with SOTA methods on VisA dataset for multi-class anomaly detection with $AUROC_{cls}/AP_{cls}/F1max_{cls}$ metrics.

Category	Non-Diffusion Method		Diffusion-based Method				
	DRAEM	UniAD	DDPM		LDM		Ours
pcb1	94.6/31.8/37.2/52.8	93.3/ 3.9/ 8.3/64.1	75.7/ 1.1/ 2.8/36.1	84.5/ 2.1/ 4.9/54.3	98.7/49.6/52.8/80.2		
pcb2	92.3/10.0/18.6/66.2	93.9/ 4.2/ 9.2/66.9	76.2/ 0.7/ 1.6/30.8	89.5/ 2.5/ 6.7/52.7	95.2/ 7.5/16.7/67.0		
pcb3	90.8/14.1/24.4/42.9	97.3/13.8/21.9/70.6	83.3/ 1.0/ 2.5/56.1	94.4/ 9.2/17.4/67.8	96.7/ 8.0/18.8/68.9		
pcb4	94.4/31.0/37.6/75.7	94.9/14.7/22.9/72.3	73.0/ 1.4/ 3.5/29.9	80.4/ 2.1/ 4.2/40.3	97.0/17.6/27.2/85.0		
macaroni1	95.0/19.1/24.1/67.0	97.4/ 3.7/ 9.7/84.0	87.4/ 0.4/ 1.0/61.2	81.6/ 0.3/ 1.3/47.3	94.1/10.2/16.7/68.5		
macaroni2	94.6/ 3.9/12.4/65.2	95.2/ 0.9/ 4.3/76.6	84.8/ 0.2/ 0.6/54.1	87.2/ 0.3/ 0.6/57.2	93.6/ 0.9/ 2.8/73.1		
capsules	97.1/27.8/33.7/62.8	88.7/ 3.0/ 7.4/43.7	77.1/ 1.1/ 2.8/34.6	75.5/ 1.1/ 2.7/34.8	97.3/10.0/21.0/77.9		
candle	82.2/10.1/19.0/65.6	98.5/17.6/27.9/91.6	76.4/ 0.4/ 1.4/34.1	85.3/ 0.9/ 1.9/46.8	97.3/12.8/22.8/89.4		
cashew	80.7/ 9.9/15.7/38.5	98.6/51.7/58.3/87.9	74.5/ 2.7/ 5.2/58.7	90.5/ 5.1/10.1/68.3	90.9/53.1/60.9/61.8		
chewinggum	91.0/62.3/63.3/40.9	98.8/54.9/56.1/81.3	74.7/ 1.4/ 2.8/37.9	84.1/ 3.1/ 6.9/52.9	94.7/11.9/25.8/59.5		
fryum	92.4/38.8/38.5/69.5	95.9/34.0/40.6/76.2	85.7/ 9.4/17.2/58.4	89.9/14.8/24.8/60.1	97.6/58.6/60.1/81.3		
pipe.fryum	91.1/38.1/39.6/61.8	98.9/50.2/57.7/91.5	87.0/ 6.9/12.9/69.6	96.4/31.0/37.2/77.6	99.4/72.7/69.9/89.9		
Mean	91.3/23.5/29.5/58.8	95.9/21.0/27.0/75.6	79.7/ 2.2/ 4.5/46.8	86.6/ 6.0/ 9.9/55.0	96.0/26.1/33.0/75.2		

Table 13: Comparison with SOTA methods on VisA dataset for multi-class anomaly localization with $AUROC_{seg}/AP_{seg}/F1max_{seg}/PRO$ metrics.

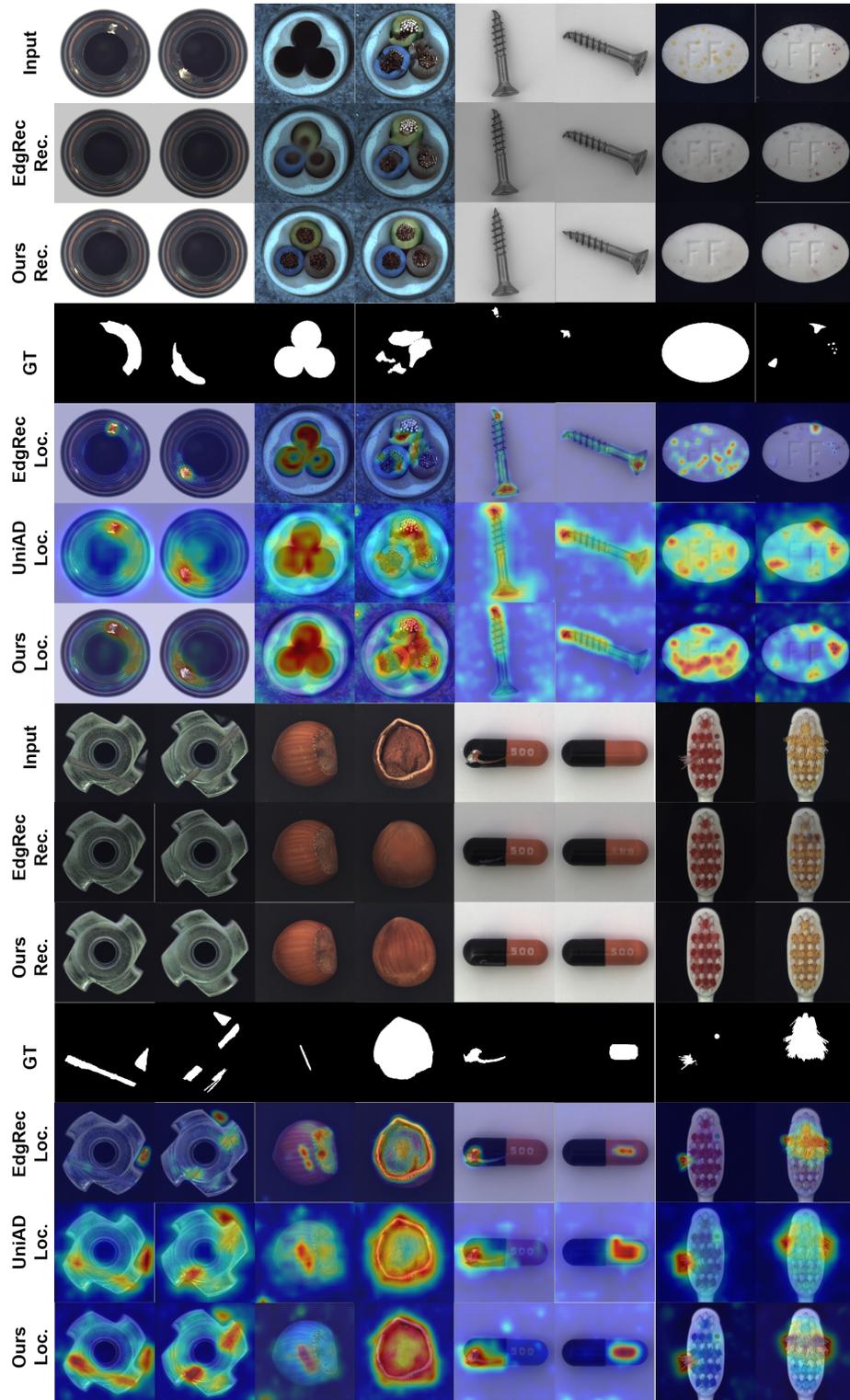


Figure 8: Qualitative comparison results for anomaly localization on MVTec-AD dataset.

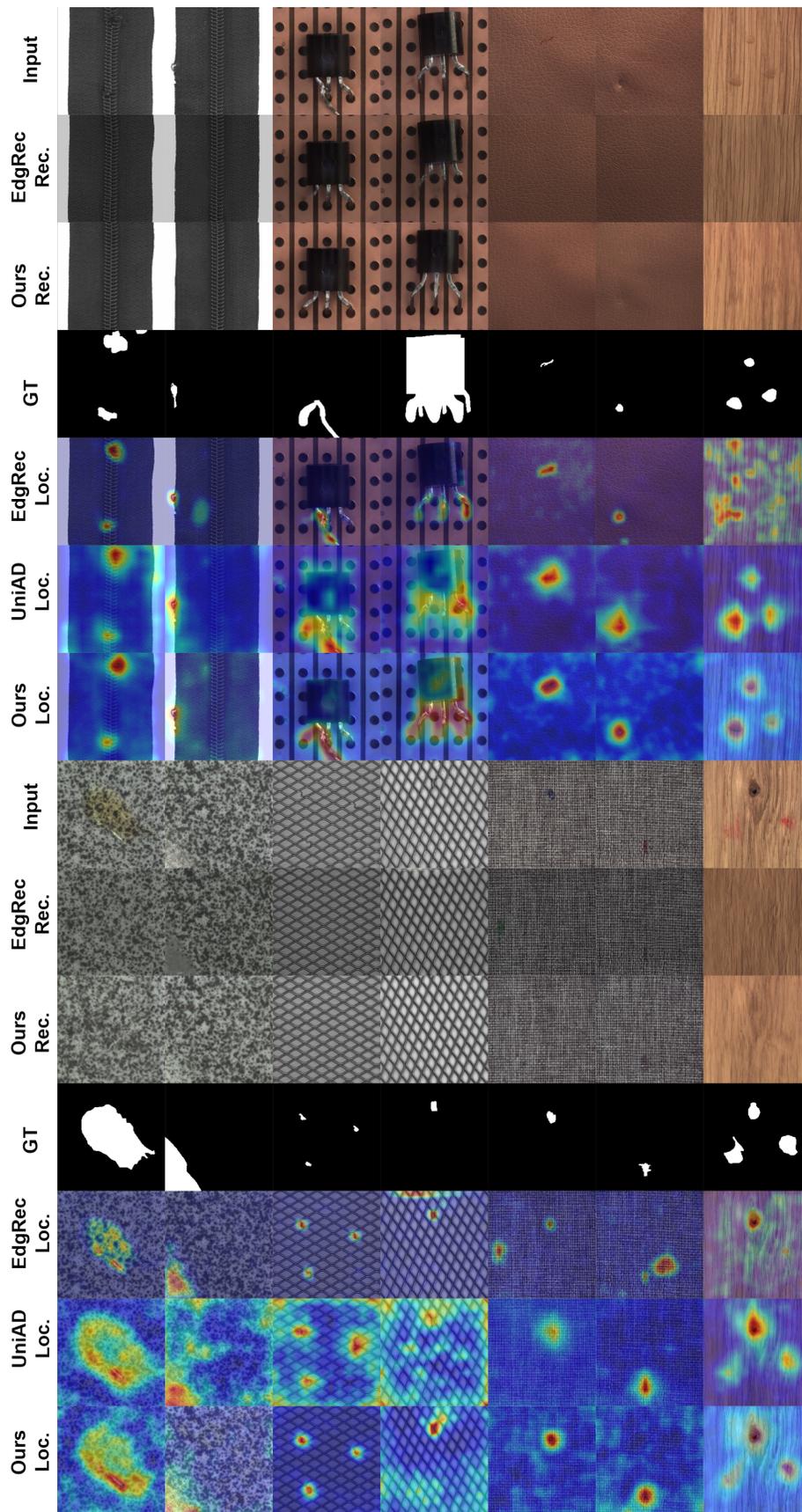


Figure 9: Qualitative comparison results for anomaly localization on MVTec-AD dataset.

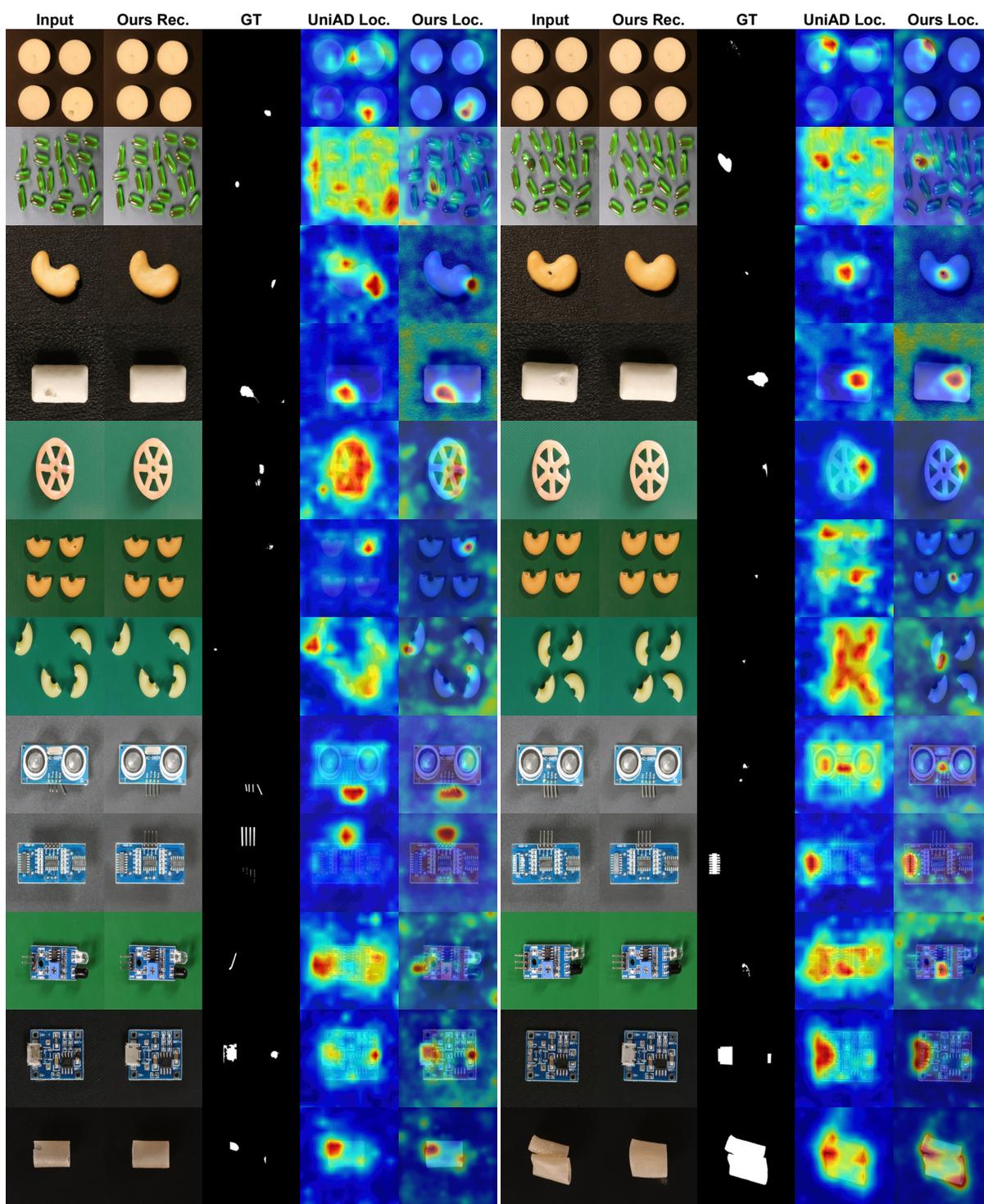


Figure 10: Qualitative comparison results for anomaly localization on VisA dataset.