

Unsupervised Anomaly Detection via Masked Diffusion Posterior Sampling

Di Wu¹, Shicai Fan^{1,2}, Xue Zhou^{1,2}, Li Yu^{1,2}, Yuzhong Deng¹, Jianxiao Zou^{1,2}, Baihong Lin^{2*}

¹School of Automation Engineering, University of Electronic Science and Technology of China(UESTC)

²Shenzhen Institute for Advanced Study, UESTC

linbaihong111@126.com

Abstract

Reconstruction-based methods have been commonly used for unsupervised anomaly detection, in which a normal image is reconstructed and compared with the given test image to detect and locate anomalies. Recently, diffusion models have shown promising applications for anomaly detection due to their powerful generative ability. However, these models lack strict mathematical support for normal image reconstruction and unexpectedly suffer from low reconstruction quality. To address these issues, this paper proposes a novel and highly-interpretable method named Masked Diffusion Posterior Sampling (MDPS). In MDPS, the problem of normal image reconstruction is mathematically modeled as multiple diffusion posterior sampling for normal images based on the devised masked noisy observation model and the diffusion-based normal image prior under Bayesian framework. Using a metric designed from pixel-level and perceptual-level perspectives, MDPS can effectively compute the difference map between each normal posterior sample and the given test image. Anomaly scores are obtained by averaging all difference maps for multiple posterior samples. Exhaustive experiments on MVTec and BTAD datasets demonstrate that MDPS can achieve state-of-the-art performance in normal image reconstruction quality as well as anomaly detection and localization.

1 Introduction

Anomaly detection (AD) is a fundamental computer vision task with widespread applications in medical diagnosis, intelligent manufacturing, autonomous driving and etc. However, due to the rarity and diversity of anomaly samples, recent studies mainly focus on unsupervised anomaly detection (UAD) [Diers and Pigorsch, 2023], in which the models only learn from normal samples but can detect anomaly data.

So far, there have existed various methods for UAD, among which reconstruction-based method is one of the earliest and

most common neural network approaches [Ruff *et al.*, 2021]. Given a test image, reconstruction-based method tries to reconstruct the corresponding normal image and compute the difference between the test image and the reconstruction to detect and localize the anomalies. Obviously, how to reconstruct a normal image is a key issue for reconstruction-based methods. Early reconstruction models include Autoencoder (AE), Generative Adversarial Network (GAN) and etc. However, AE can easily suffer from “identical shortcut” [You *et al.*, 2022a] and blurry reconstruction [Baur *et al.*, 2021], i.e., they reconstruct both normal and anomalous images and the results tend to be more blurry. Although GAN can alleviate the problems of AE, mode collapse and training instability make GAN challenging for UAD [Xia *et al.*, 2022].

Recently, Diffusion Models (DMs) have attracted most researchers’ attention with their powerful image generation ability [Ho *et al.*, 2020]. Compared with previous generative models, DMs can effectively record image priors and generate various realistic images after simple training based on Gaussian denoising. Thus, they show promising applications for UAD to alleviate the problems of inferior reconstruction quality or insufficient coverage of the normal image distribution. However, two key problems arise when introducing DMs for UAD: First, although different DM-based methods are proposed for UAD, they are lack of strict mathematical theories or interpretability to ensure that the anomaly region of a test image can be reconstructed as the normal one. Second, most DMs unexpectedly suffer from low reconstruction quality, especially for the normal region of a test image, since Gaussian noise in DMs will destroy the original normal texture. This problem can easily result in misjudgment of anomalous pixels in the normal region. Thus, it requires further studies for DMs to maintain the texture in the normal region of a test image after reconstruction.

In this paper, we propose a novel Masked Diffusion Posterior Sampling (MDPS) method for UAD under Bayesian framework, which has high interpretability supported by relatively-strict mathematics. In our method, we firstly propose a Masked Noisy Observation Model which regards a test image as a masked noisy observation of a normal image to protect the normal region of a test image and enhance reconstruction quality. Then, we take a Denoising Diffusion Implicit Model (DDIM) trained by normal samples as image prior, and model the problem of normal image reconstruc-

*B. Lin is the corresponding author. Source code will be available at <https://github.com/KevinBHLin/>.

tion as multiple diffusion posterior samplings for normal images based on the devised observation model and the DDIM-based prior. Third, using a metric designed from pixel-level and perceptual-level perspectives, we compute the difference maps between multiple normal posterior samples and the test image respectively, and average all difference maps to accurately obtain anomaly scores. Exhaustive experiments on MVTEC and BTAD datasets prove that the proposed MDPS achieves excellent reconstruction quality and high accuracy of anomaly detection and localization compared with other state-of-the-art reconstruction-based methods.

2 Related Work

Recent reconstruction-based UAD methods can be roughly divided into three categories, including AE-based methods, GAN-based methods, and DM-based methods.

AE-based methods. Early reconstruction-based UAD methods commonly adopt AE due to their simple architectures and easy-to-implement training processes. Unfortunately, AEs can easily suffer from the problems of “identical shortcut” and blurry reconstruction. For the first problem, since Vision Transformer (ViT) can prevent “identical shortcut” [You *et al.*, 2022a], recent studies try to design ViT-based AEs, e.g., [Mishra *et al.*, 2021], [You *et al.*, 2022b], [You *et al.*, 2022a] and etc. For the second problem, various methods have been proposed: [Bergmann *et al.*, 2019] designs a loss function of structure similarity to improve reconstruction quality; [Liu *et al.*, 2020] and [Dehaene *et al.*, 2020] introduce Variational AE to achieve better reconstruction results; [Zavrtanik *et al.*, 2021] and [Sun *et al.*, 2023] propose a self-supervised model to reduce the influence of blurry reconstruction. Nevertheless, the above schemes cannot totally avoid blurry reconstruction, thus, leading to performance bottleneck for UAD.

GAN-based methods. To overcome the drawbacks of AEs, [Schlegl *et al.*, 2017] firstly proposed a GAN-based method named AnoGAN. Since then, various variants are proposed, e.g., [Akçay *et al.*, 2019], [Schlegl *et al.*, 2019] and etc. Although GANs can empirically generate high definition results, two problems make GANs challenging for UAD [Xia *et al.*, 2022]: First, GAN training is highly unstable to converge, sometimes bringing in meaningless reconstruction. Since the training process of AE is more stable, AE and GAN are usually combined to alleviate training instability, e.g., [Tang *et al.*, 2020], [Contreras-Cruz *et al.*, 2023]. Second, due to the vanish of discriminator gradient, it is hard to ensure sufficient coverage of the normal distribution during GAN training process. This phenomenon is called mode collapse, which can easily lead to reconstruction results with only a few modes. To alleviate this problem, self-supervised models are proposed based on simulated anomaly data to guide GAN training, e.g., [Song *et al.*, 2022], [Liu *et al.*, 2023] and etc.

DM-based methods. Recently, DMs have been popular for academic research due to their potentially-powerful generative ability. Typical DMs include DDPM [Ho *et al.*, 2020], DDIM [Song *et al.*, 2021a], SDE [Song *et al.*, 2021b], etc. Compared with AEs and GANs, DMs have desirable properties, such as distribution coverage, a stationary training

objective and easy scalability [Dhariwal and Nichol, 2021]. To generate desirable images based on DMs for certain requirements, conditional DMs are studied, e.g., [Dhariwal and Nichol, 2021] proposes a classifier-guided DM to generate images given a certain class label; [Chung *et al.*, 2023] extends diffusion solvers to efficiently handle general noisy inverse problems via approximation of posterior sampling. The above works show promising applications for UAD to reconstruct high definition normal images.

So far, there have been a few DM-based AD methods. [Wolleb *et al.*, 2022] and [Pinaya *et al.*, 2022] firstly introduce DMs for weakly supervised AD and UAD respectively in medical diagnosis. [Wyatt *et al.*, 2022] proposes a DM-based UAD method named AnoDDPM, which uses simple noise to improve normal image reconstruction quality for brain MRI. [Gonzalez-Jimenez *et al.*, 2023] uses score-based models to produce a gradient map highlighting anomaly areas. [Bercea *et al.*, 2023] proposes a DM-based model named AutoDDPM which consists of three stages, i.e, mask, stitch, and re-sampling. [Lu *et al.*, 2023] models the problem of normal image reconstruction as a DM-based denoising process. The above methods have enhanced the accuracy of DM-based AD from different aspects. However, these DM-based methods are lack of strictly-mathematical support, and unexpectedly suffer from low normal image reconstruction quality, especially for the normal region of a test image.

3 Background

In this section, we briefly review SDE-based diffusion models [Song *et al.*, 2021b] and DDIM [Song *et al.*, 2021a].

SDE-based diffusion models. For a diffusion process $\{\mathbf{x}(t)\}_{t=0}^T$ indexed by a continuous time variable $t \in [0, T]$, let $\mathbf{x}(0) \sim p_0(\mathbf{x})$ and $\mathbf{x}(T) \sim p_T(\mathbf{x})$, in which $p_0(\mathbf{x})$ and $p_T(\mathbf{x})$ denotes the data distribution of interest and a known spherical Gaussian distribution respectively. Then, the forward noising process $\mathbf{x}(0) \rightarrow \mathbf{x}(T)$ can be modeled as the following Itô stochastic differential equation (SDE):

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + g(t)d\boldsymbol{\omega} \quad (1)$$

where $\mathbf{f}(\cdot, t)$ and $g(\cdot)$ are the drift and diffusion coefficient of $\mathbf{x}(t)$ respectively, $\boldsymbol{\omega}$ is a standard Wiener process. The corresponding reverse process $\mathbf{x}(T) \rightarrow \mathbf{x}(0)$ is given by

$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x})]dt + \mathbf{g}(t)d\hat{\boldsymbol{\omega}}, \quad (2)$$

where $\hat{\boldsymbol{\omega}}$ is a standard Wiener process running backwards, and dt is an infinitesimal negative time step. $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ is the score function of each marginal distribution which can be approximately estimated by training a time-dependent score-based model $s_\theta(\mathbf{x}, t)$, i.e., $s_\theta(\mathbf{x}, t) \simeq \nabla_{\mathbf{x}} \log p_t(\mathbf{x})$.

DDIMs. DDIMs can be regarded as an acceleration version of DDPMs [Ho *et al.*, 2020] with the same training procedure. However, different from DDPMs, the forward noising process of DDIM is defined as a non-Markov process:

$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1 - \alpha_t}\boldsymbol{\epsilon}, \boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I}) \quad (3)$$

where $\{\alpha_t \in (0, 1) | t \in [0, T]\}$ is a decreasing sequence decided by a predetermined schedule, and $\boldsymbol{\epsilon}$ is white noise

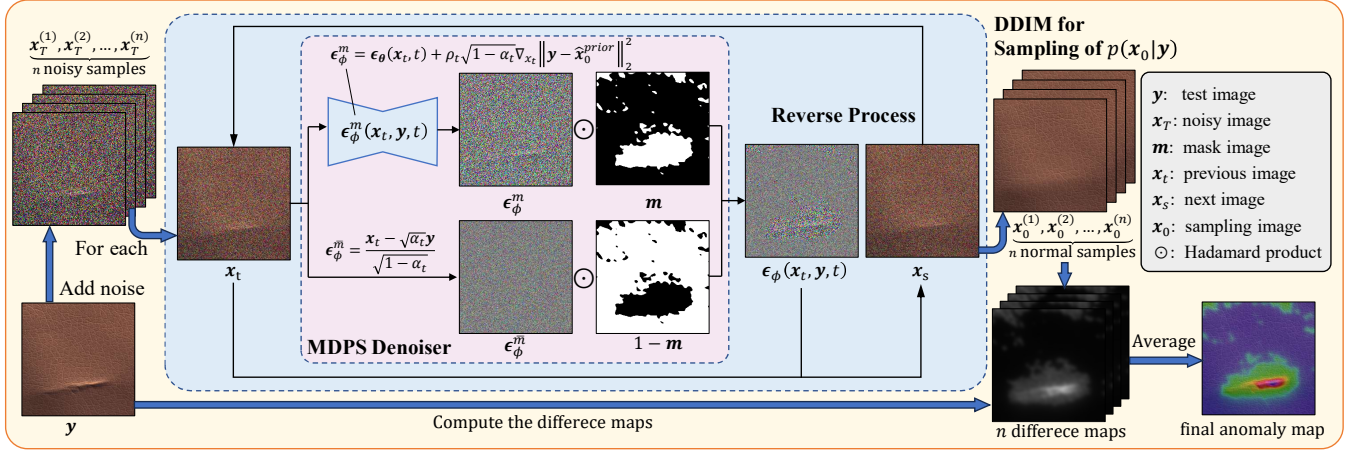


Figure 1: Overview of MDPS. The MDPS denoiser shown in the pink box is designed partially based on the denoiser $\epsilon(x_t, t)$ of DDIM for sampling of $p(x_0)$. Based on MDPS, we can obtain n normal posterior samples from n noisy versions of the test image \mathbf{y} respectively. Then, the final anomaly map is obtained by averaging n difference maps computed from n normal posterior samples and the test image \mathbf{y} .

drawn from a standard normal distribution. The accelerated reverse process of DDIM is described as:

$$\begin{aligned} \mathbf{x}_s &= \sqrt{\frac{\alpha_s}{\alpha_t}} (\mathbf{x}_t - \sqrt{1 - \alpha_t} \epsilon_\theta(\mathbf{x}_t, t)) \\ &\quad + \sqrt{1 - \alpha_s - \sigma_t^2} \epsilon_\theta(\mathbf{x}_t, t) + \sigma_t \epsilon_t, \quad \epsilon_t \sim \mathcal{N}(0, \mathbf{I}) \end{aligned} \quad (4)$$

where $\sigma_t = \sqrt{(1 - \alpha_s)/(1 - \alpha_t)} \sqrt{1 - \alpha_t/\alpha_s}$, $s < t$, and $\epsilon_\theta(\mathbf{x}_t, t)$ is a U-net denoiser which estimates the white noise ϵ from \mathbf{x}_t , and can be trained by minimizing the following objective:

$$L(\theta) := \mathbb{E}_{t, \mathbf{x}_0, \epsilon} [\|\epsilon - \epsilon_\theta(\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \epsilon, t)\|^2] \quad (5)$$

[Song *et al.*, 2021b] has pointed out that DDIM is a discrete form of SDE when $\mathbf{f}(\mathbf{x}, t) = -\mathbf{x}_t \beta(t)/2$, $g(t) = \sqrt{\beta(t)}$, and $\alpha_t = \prod_{s=0}^{t-1} (1 - \beta(s))$ in Eqn.(1). Thus,

$$\begin{aligned} \epsilon_\theta(\mathbf{x}_t, t) &= -\sqrt{1 - \alpha_t} \mathbf{s}_\theta(\mathbf{x}_t, t) \\ &\approx -\sqrt{1 - \alpha_t} \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) \end{aligned} \quad (6)$$

4 Masked Diffusion Posterior Sampling

In this section, we firstly establish a masked noisy observation model for normal image reconstruction. Then, based on the observation model, we propose a posterior sampling method for normal images using DDIMs. Thirdly, we design image-wise and pixel-wise anomaly scores and propose a mask generation scheme based on the designed anomaly scoring for the proposed observation model to enhance the reconstruction quality of posterior samples.

4.1 Masked Noisy Observation Model

Let \mathbf{y} and \mathbf{x}_0 denote an anomaly image and the corresponding normal image respectively. Let \mathbf{m} denote a mask image with the same size of \mathbf{y} in which pixel values are set to 0 in the normal regions of \mathbf{y} and set to 1 otherwise. Then, we establish a masked noisy observation model to describe the relationship between \mathbf{y} and \mathbf{x}_0 as the following:

$$\mathbf{y} = (1 - \mathbf{m}) \odot \mathbf{x}_0 + \mathbf{m} \odot (\mathbf{x}_0 + \mathbf{n}) \quad (7)$$

where \odot denotes Hadamard product, \mathbf{n} denotes zero-mean isotropic Gaussian noise and $\mathbf{n} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$. Obviously, the above modeling assumes that \mathbf{y} is a noisy observation of \mathbf{x}_0 in the anomaly region whereas \mathbf{y} shares the same pixel values as that of \mathbf{x}_0 in the normal region, i.e.,

$$(1 - \mathbf{m}) \odot \mathbf{y} = (1 - \mathbf{m}) \odot \mathbf{x}_0 \quad (8a)$$

$$\mathbf{m} \odot \mathbf{y} = \mathbf{m} \odot (\mathbf{x}_0 + \mathbf{n}). \quad (8b)$$

As \mathbf{y} would contain unknown types of anomaly, the assumption of Gaussian noise for anomaly is reasonable. Furthermore, the mask $1 - \mathbf{m}$ can guarantee the reconstruction quality of normal image \mathbf{x}_0 , since it maintains pixel values in the normal region of \mathbf{y} for reconstruction of \mathbf{x}_0 . Then, based on Eqn.(7), we can obtain the distribution of \mathbf{y} given \mathbf{x}_0 for the anomaly region indicated by \mathbf{m} as the following:

$$p(\mathbf{y}|\mathbf{x}_0) \sim \mathcal{N}(\mathbf{x}_0, \sigma^2 \mathbf{I}) \quad (9)$$

4.2 Posterior Sampling for Normal Images

To obtain multiple normal samples of \mathbf{x}_0 given \mathbf{y} , we propose to conduct posterior sampling for $p(\mathbf{x}_0|\mathbf{y})$. According to Bayes' Theorem, $p(\mathbf{x}_0|\mathbf{y})$ is related to the normal image prior $p(\mathbf{x}_0)$ and the observation distribution $p(\mathbf{y}|\mathbf{x}_0)$ given in Eqn.(9). In our model, we introduce a DDIM trained by normal samples to model the normal image prior $p(\mathbf{x}_0)$, since this model has been widely adopted for its powerful image generation ability in recent years. Then, following [Chung *et al.*, 2023], we regard the posterior sampling process as a diffusion process based on the introduced DDIM for $p(\mathbf{x}_0)$, and obtain multiple normal samples drawn from $p(\mathbf{x}_0|\mathbf{y})$. The details are shown as the following.

Specifically, let $\epsilon_\theta(\mathbf{x}_t, t)$ and $\epsilon_\phi(\mathbf{x}_t, \mathbf{y}, t)$ denote the denoisers of two DDIMs for sampling of $p(\mathbf{x}_0)$ and $p(\mathbf{x}_0|\mathbf{y})$ respectively. Since $\epsilon_\theta(\mathbf{x}_t, t)$ is known and trained by normal images, we discuss the design of $\epsilon_\phi(\mathbf{x}_t, \mathbf{y}, t)$. According to Eqn.(3) and Eqn.(8a), we can accurately estimate ϵ as:

$$\epsilon_\phi^{\overline{\mathbf{m}}}(\mathbf{x}_t, \mathbf{y}, t) = \frac{\mathbf{x}_t - \sqrt{\alpha_t} \mathbf{x}_0}{\sqrt{1 - \alpha_t}} = \frac{\mathbf{x}_t - \sqrt{\alpha_t} \mathbf{y}}{\sqrt{1 - \alpha_t}} \quad (10)$$

where $\epsilon_\phi^{\bar{m}}(\mathbf{x}_t, \mathbf{y}, t)$ denotes the values of $\epsilon_\phi(\mathbf{x}_t, \mathbf{y}, t)$ in the normal region indicated by $1 - \mathbf{m}$.

However, for the anomaly region indicated by \mathbf{m} , according to Bayes' Theorem, Eqn.(6) and Eqn.(8b),

$$\begin{aligned} \epsilon_\phi^{\mathbf{m}}(\mathbf{x}_t, \mathbf{y}, t) &\simeq -\sqrt{1 - \alpha_t} \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | \mathbf{y}) \\ &= -\sqrt{1 - \alpha_t} (\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)) \quad (11) \\ &\simeq \epsilon_\theta(\mathbf{x}_t, t) - \sqrt{1 - \alpha_t} \nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t) \end{aligned}$$

where $\epsilon_\phi^{\mathbf{m}}(\mathbf{x}_t, \mathbf{y}, t)$ denotes the result of $\epsilon_\phi(\mathbf{x}_t, \mathbf{y}, t)$ for the anomaly region indicated by \mathbf{m} . Obviously, $\epsilon_\phi^{\mathbf{m}}(\mathbf{x}_t, \mathbf{y}, t)$ can be replaced by the above equation base on $\epsilon_\theta(\mathbf{x}_t, t)$ and $\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)$. Since $\epsilon_\theta(\mathbf{x}_t, t)$ is a known denoiser trained by normal samples, we focus on the modeling of $p(\mathbf{y} | \mathbf{x}_t)$.

As has been proven in [Chung *et al.*, 2023],

$$p(\mathbf{y} | \mathbf{x}_t) = \int p(\mathbf{y} | \mathbf{x}_0) p(\mathbf{x}_0 | \mathbf{x}_t) d\mathbf{x}_0 \quad (12)$$

where $p(\mathbf{y} | \mathbf{x}_0)$ has been given in Eqn.(9), and $p(\mathbf{x}_0 | \mathbf{x}_t)$ is determined by the reverse process of DDIMs for sampling of $p(\mathbf{x}_0)$. Based on Eqn.(3) and Eqn.(4),

$$\hat{\mathbf{x}}_0^{\text{prior}} = \frac{1}{\sqrt{\alpha_t}} (\mathbf{x}_t - \sqrt{1 - \alpha_t} \epsilon_\theta(\mathbf{x}_t, t)). \quad (13)$$

Apparently, in the reverse process of DDIMs, ϵ of Eqn.(3) turns to be an unknown fixed constant, and can be estimated by $\epsilon_\theta(\mathbf{x}_t, t)$. The estimation errors between ϵ and $\epsilon_\theta(\mathbf{x}_t, t)$ are reduced by minimizing the objective function of Eqn.(5) during the training process of DDIMs for sampling of $p(\mathbf{x}_0)$. Thus, we can assume that $\epsilon \sim \mathcal{N}(\epsilon_\theta(\mathbf{x}_t, t), \lambda_t^{-1} \mathbf{I})$ in which λ_t indicates the estimation precision, and obtain

$$p(\mathbf{x}_0 | \mathbf{x}_t) \sim \mathcal{N}(\hat{\mathbf{x}}_0^{\text{prior}}, \frac{1 - \alpha_t}{\lambda_t \alpha_t} \mathbf{I}) \quad (14)$$

Then, for the anomaly region indicated by \mathbf{m} , $p(\mathbf{y} | \mathbf{x}_t)$ can be calculated based on Eqn.(9), Eqn.(12) and Eqn.(14) as:

$$p(\mathbf{y} | \mathbf{x}_t) \sim \mathcal{N}(\hat{\mathbf{x}}_0^{\text{prior}}, \rho_t^{-1} \mathbf{I}), \quad \rho_t^{-1} = \frac{1 - \alpha_t}{\lambda_t \alpha_t} + \sigma^2 \quad (15)$$

For simplicity, we let $\rho_t = \rho$, ρ is an adjustable guidance scale. $\epsilon_\phi^{\mathbf{m}}(\mathbf{x}_t, \mathbf{y}, t)$ can be obtained based on Eqn.(15) as:

$$\epsilon_\phi^{\mathbf{m}}(\mathbf{x}_t, \mathbf{y}, t) = \epsilon_\theta(\mathbf{x}_t, t) + \rho \sqrt{1 - \alpha_t} \nabla_{\mathbf{x}_t} \|\mathbf{y} - \hat{\mathbf{x}}_0^{\text{prior}}\|_2^2 \quad (16)$$

Combining Eqn.(16) and Eqn.(10), we can obtain the denoiser $\epsilon_\phi(\mathbf{x}_t, \mathbf{y}, t)$ of DDIMs for sampling of $p(\mathbf{x}_0 | \mathbf{y})$ as:

$$\begin{aligned} \epsilon_\phi(\mathbf{x}_t, \mathbf{y}, t) &= (1 - \mathbf{m}) \odot \frac{\mathbf{x}_t - \sqrt{\alpha_t} \mathbf{y}}{\sqrt{1 - \alpha_t}} \\ &+ \mathbf{m} \odot (\epsilon_\theta(\mathbf{x}_t, t) + \rho \sqrt{1 - \alpha_t} \nabla_{\mathbf{x}_t} \|\mathbf{y} - \hat{\mathbf{x}}_0^{\text{prior}}\|_2^2) \end{aligned} \quad (17)$$

4.3 Anomaly Scoring and Mask Generation

To detect and locate anomaly, we have to predict image-wise and pixel-wise anomaly scores respectively. For prediction of pixel-wise anomaly score, we have to calculate the difference between the test image and the corresponding normal reconstruction result. Traditionally, pixel-level metrics such as \mathcal{L}_1 -norm, and MSE are adopted to measure the differences of two

Algorithm 1 Masked Diffusion Posterior Sampling

Input: test image \mathbf{y} , DDIM denoiser $\epsilon_\theta(\mathbf{x}_t, t)$, mask image \mathbf{m} , guidance scale ρ , noise level T , sampling times N

Output: normal image \mathbf{x}_0

```

1:  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ .
2:  $\mathbf{x}_T = \sqrt{\alpha_T} \mathbf{y} + \sqrt{1 - \alpha_T} \epsilon$ .
3: for all  $n$  from  $N$  to 1 do
4:    $t = \frac{T^n}{N}, s = \frac{T^{(n-1)}}{N}$ 
5:    $\hat{\mathbf{x}}_0^{\text{prior}} = \frac{1}{\sqrt{\alpha_t}} (\mathbf{x}_t - \sqrt{1 - \alpha_t} \epsilon_\theta(\mathbf{x}_t, t))$ 
6:    $\epsilon_\phi^{\mathbf{m}} = \epsilon_\theta(\mathbf{x}_t, t) + \rho \sqrt{1 - \alpha_t} \nabla_{\mathbf{x}_t} \|\mathbf{y} - \hat{\mathbf{x}}_0^{\text{prior}}\|_2^2$ 
7:    $\epsilon_\phi^{\bar{m}} = \frac{\mathbf{x}_t - \sqrt{\alpha_t} \mathbf{y}}{\sqrt{1 - \alpha_t}}$ 
8:    $\epsilon_\phi = (1 - \mathbf{m}) \odot \epsilon_\phi^{\bar{m}} + \mathbf{m} \odot \epsilon_\phi^{\mathbf{m}}$ 
9:    $\hat{\mathbf{x}}_0 = \frac{1}{\sqrt{\alpha_t}} (\mathbf{x}_t - \sqrt{1 - \alpha_t} \epsilon_\phi)$ 
10:   $\epsilon_t \sim \mathcal{N}(0, \mathbf{I})$ 
11:   $\mathbf{x}_s = \sqrt{\alpha_s} \hat{\mathbf{x}}_0 + \sqrt{1 - \alpha_s - \sigma_t^2} \epsilon_\phi + \sigma_t \epsilon_t$ 
12: end for
13: return  $\mathbf{x}_0$ 

```

images. Recently, [Zhang *et al.*, 2018] proposed a perceptual-level metric named LPIPS based on feature maps extracted by a pretrained convolutional neural network. Since LPIPS can better match human perceptual similarity judgments, it has been widely adopted for reconstruction-based AD method, e.g., [Defard *et al.*, 2021], [Roth *et al.*, 2022] and etc. Following recent works, the difference between \mathbf{y} and \mathbf{x}_0 at the k -th spatial position is calculated as:

$$\begin{aligned} [\mathcal{D}(\mathbf{x}_0, \mathbf{y})]_k &= \underbrace{\eta \|\mathbf{y}_k - (\mathbf{x}_0)_k\|_1}_{\text{Pixel-level Metric}} \\ &+ \underbrace{\sum_{i \in \mathcal{J}} \left(1 - \frac{(\mathcal{F}_i^{(k)}(\mathbf{x}_0))^T \mathcal{F}_i^{(k)}(\mathbf{y})}{\|\mathcal{F}_i^{(k)}(\mathbf{x}_0)\| \|\mathcal{F}_i^{(k)}(\mathbf{y})\|} \right)}_{\text{Perceptual-level Metric}} \quad (18) \end{aligned}$$

in which \mathcal{L}_1 -norm and LPIPS are incorporated to measure the differences of two images from pixel-level and perceptual-level perspectives, $\mathcal{F}_i^{(k)}$ denotes the i -th stage output feature of ResNet-101 proposed by [Zagoruyko and Komodakis, 2016] pretrained on ImageNet at the k -th spatial position, and η is a hyperparameter to balance the weight between the pixel-level and perceptual-level metrics. \mathcal{J} denotes the set of different stages of ResNet. In our paper, \mathcal{J} is set to $\{1, 2, 3\}$. Since our model will generate multiple normal samples given a single test image, the pixel-wise anomaly score map for \mathbf{y} is defined based on Eqn.(18) as:

$$\bar{\mathcal{D}} = \frac{1}{N_s} \sum_{j=1}^{N_s} \mathcal{D}(\mathbf{x}_0^{(j)}, \mathbf{y}) \quad (19)$$

where $\mathbf{x}_0^{(j)}$ denotes the j -th reconstructed normal sample, and N_s denotes the number of normal samples.

Then, based on Eqn.(19), we further design the image-wise anomaly score for \mathbf{y} as the average of the largest S pixel-wise anomaly scores in $\bar{\mathcal{D}}$ to mitigate false positives caused by image noise. In our paper, S is set to 500.

Table 1: Anomaly detection and localization performance on MVTec. (Image-AUROC %, Pixel-AUROC %)

Method	AE-based Methods		GAN-based Methods		DM-based Methods				
	DRÆM	UniAD	AnoGAN	AnoSeg	AnoDDPM	AutoDDPM	RAN	MDPS	MDPS
								($N_s = 1$)	($N_s = 16$)
Carpet	(97.0, 95.5)	(99.8, 98.5)	(33.7, 54.0)	(96.0, 99.0)	(54.6, 63.8)	(85.9, 86.0)	(99.9, 98.9)	(98.7, 93.4)	(99.6, 94.4)
Grid	(99.9, 99.7)	(98.2, 96.5)	(87.1, 58.0)	(99.0, 99.0)	(96.0, 82.3)	(100, 97.5)	(99.7, 99.1)	(100, 99.4)	(100, 99.4)
Leather	(100, 98.6)	(100, 98.8)	(45.1, 64.0)	(99.0, 98.0)	(96.9, 85.6)	(90.0, 91.5)	(100, 99.5)	(100, 99.4)	(100, 99.5)
Tile	(99.6, 99.2)	(99.3, 91.8)	(40.1, 50.0)	(98.0, 98.0)	(96.5, 76.6)	(86.8, 74.0)	(98.0, 92.1)	(99.8, 95.4)	(100, 96.4)
Wood	(99.1, 96.4)	(98.6, 93.2)	(56.7, 62.0)	(99.0, 98.0)	(87.2, 68.5)	(98.5, 80.5)	(98.1, 94.5)	(99.0, 94.1)	(99.1, 95.7)
Bottle	(99.2, 99.1)	(99.7, 98.1)	(80.0, 86.0)	(98.0, 99.0)	(87.8, 68.9)	(99.0, 97.4)	(99.3, 97.7)	(100, 98.7)	(100, 98.6)
Cable	(91.8, 94.7)	(95.2, 97.3)	(47.7, 78.0)	(98.0, 99.0)	(71.6, 64.1)	(84.7, 88.3)	(91.2, 95.6)	(97.1, 95.9)	(98.3, 95.8)
Capsule	(98.5, 94.3)	(86.9, 98.5)	(44.2, 84.0)	(84.0, 90.0)	(60.1, 72.5)	(61.8, 92.8)	(84.1, 97.5)	(91.7, 93.2)	(91.4, 93.6)
Hazelnut	(100, 99.7)	(99.8, 98.1)	(25.9, 87.0)	(98.0, 99.0)	(69.5, 75.5)	(96.5, 92.9)	(97.9, 97.3)	(99.6, 98.6)	(99.8, 98.6)
Metalnut	(98.9, 99.5)	(99.2, 94.8)	(28.4, 76.0)	(95.0, 99.0)	(55.5, 76.0)	(93.5, 94.8)	(99.2, 96.8)	(100, 97.3)	(99.9, 97.7)
Pill	(98.9, 97.6)	(93.7, 95.0)	(71.1, 87.0)	(87.0, 94.0)	(75.7, 73.6)	(59.3, 92.2)	(64.7, 92.5)	(96.2, 99.0)	(96.8, 99.2)
Screw	(93.9, 97.6)	(87.5, 98.3)	(10.0, 80.0)	(97.0, 91.0)	(64.7, 79.1)	(77.5, 93.0)	(89.9, 99.0)	(93.1, 98.7)	(96.7, 98.9)
Toothbrush	(100, 98.1)	(94.2, 98.4)	(43.9, 90.0)	(99.0, 96.0)	(57.2, 87.8)	(92.5, 96.9)	(96.9, 98.6)	(100, 98.7)	(100, 98.8)
Transistor	(93.1, 90.9)	(99.8, 97.9)	(69.2, 80.0)	(96.0, 96.0)	(70.8, 63.0)	(80.1, 77.8)	(92.3, 93.1)	(99.9, 94.1)	(100, 94.7)
Zipper	(100, 98.8)	(95.8, 96.8)	(71.5, 78.0)	(99.0, 98.0)	(92.0, 66.5)	(96.3, 89.0)	(85.5, 97.6)	(100, 98.5)	(100, 98.5)
Average	(98.0, 97.3)	(96.5, 96.8)	(50.3, 74.3)	(96.1, 96.9)	(75.7, 73.6)	(86.8, 89.6)	(93.1, 96.7)	(98.4, 97.0)	(98.8, 97.3)

Using the above image-wise and pixel-wise anomaly scores, we design a mask image generation scheme for our observation model of Eqn.(7). To obtain an accurate mask m , we firstly set m to 1, i.e., all pixels in the test image y are regarded as potential anomaly pixels, and run a group of posterior sampling for normal images. Then, we obtain the score map \bar{D} for y based on Eqn.(19), and estimate m as:

$$m = (\bar{D} > T_{th}), T_{th} = \min \bar{D} + \lambda(\max \bar{D} - \min \bar{D}) \quad (20)$$

where λ is a hyperparameter, and T_{th} is a threshold for anomaly score determined by λ . Using the above mask generation scheme, we re-run another group of posterior sampling for normal images and calculate anomaly scores for each image and each pixel respectively. The proposed MDPS is summarized in **Algorithm 1**.

5 Experiments

In this section, we compare our MDPS with other UAD methods, and conduct ablation studies to validate the designs.

Datasets. We conduct all experiments on the MVTec and BTAD Datasets. The MVTec dataset is an industrial AD benchmark [Bergmann *et al.*, 2021], which contains 15 categories (5 textural categories and 10 object categories) with about 200 normal samples and 100 anomaly samples for each class. It provides various types of anomaly with pixel-level segmentation ground truths such as scratches, cracks, color, and missing components, posing a great challenge to AD. In our experiment, following [Roth *et al.*, 2022], we resize all images to 256×256 and center crop the images to 224×224 for MVTec. The BTAD dataset contains approximately 2500 real-world industrial images of three products [Mishra *et al.*, 2021], which is more challenging for AD. Since some anomalies are located in the edge regions of the image in BTAD, we only resize all images to 256×256 without center cropping.

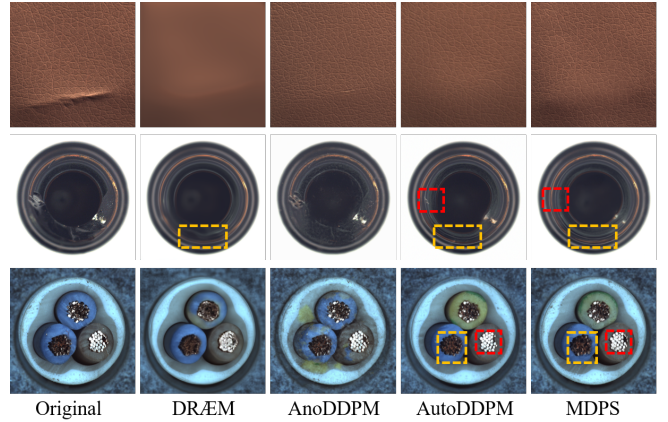


Figure 2: Comparisons of the reconstruction results on MVTec. Noting the area in dotted boxes, MDPS gives the best reconstruction results.

Evaluation Metrics. To evaluate the results of all AD comparison methods, we employ the metrics of *Area Under the Receiver Operating characteristic Curve* (AUROC). Specifically, we adopt Image-AUROC to evaluate the accuracy of anomaly detection, and adopt Pixel-AUROC to evaluate the accuracy of anomaly localization.

Implementation details. We adopt the U-net architecture proposed by [Dhariwal and Nichol, 2021] to implement the denoiser $\epsilon_{\theta}(x_t, t)$ of DDIM for sampling of $p(x_0)$. For each category of normal samples in MVTec/BTAD, we train a UNet denoiser $\epsilon_{\theta}(x_t, t)$ separately within 2000 epochs using an Adam optimizer (learning rate: $1e-4$, weight decay: $5e-2$) based on a single GeForce RTX 3090 GPU. In the training process, the batchsize is set to 8, and the timestep of DDIM is set to be 1000. After training, we utilize the trained denoiser for the proposed MDPS, and let $T = 200$, $N = 10$, $\rho = 100$. In Section 5.2, we will further discuss the selection

Table 2: Anomaly detection and localization performance on BTAD. (Image-AUROC %, Pixel-AUROC %).

Methods	classes			Average
	01	02	03	
DRAEM	(98.5,91.5)	(68.6,73.4)	(99.8,96.3)	(89.0,87.1)
VT-ADL	(97.6, 99.0)	(71.0, 94.0)	(82.6, 77.0)	(83.7, 90.0)
AnoDDPM	(71.0,62.3)	(60.1,60.4)	(52.0, 53.3)	(61.0,58.7)
AutoDDPM	(96.1,67.5)	(76.7,59.7)	(99.3, 74.3)	(90.7,67.2)
PatchCore	(90.9, 95.5)	(79.3, 94.7)	(99.8, 99.3)	(90.0, 96.5)
PaDiM	(99.8, 97.0)	(82.0, 96.0)	(99.4, 98.8)	(93.7, 97.3)
PyramidFlow	(100 , 97.4)	(88.2, 97.6)	(99.3, 98.1)	(95.8, 97.7)
MDPS($N_s = 1$)	(100 , 98.3)	(99.9 , 95.1)	(100 , 99.4)	(99.9 , 97.6)
MDPS($N_s = 16$)	(100 , 98.4)	(95.2, 95.3)	(100 , 99.4)	(98.4, 97.7)

Table 3: Comparison on time consumption for per image.

Method	Patchcore	DRÆM	AutoDDPM	MDPS($N_s = 1$)
Times(s)	0.17-0.6	0.13	33.5	0.5

of hyperparameters for MDPS in details.

5.1 Comparison with State-of-the-art

MVTec. We compare MDPS with several representative reconstruction-based methods on the MVTec dataset, including DRÆM [Zavrtanik *et al.*, 2021], UniAD [You *et al.*, 2022a], AnoGAN [Schlegl *et al.*, 2017], AnoSeg [Song *et al.*, 2022], AnoDDPM [Wyatt *et al.*, 2022], AutoDDPM [Bercea *et al.*, 2023], and RAN [Lu *et al.*, 2023]. The results are shown in Table 1 and Figure 2.

From Table 1 and Figure 2, we can find that MDPS achieves the best performance in reconstruction quality as well as anomaly detection and localization especially when $N_s = 16$. Specifically, the Image-AUROC of MDPS with $N_s = 16$ outperforms the second best comparison method DRÆM by 0.8%, and outperforms the second best DM-based method RAN by 5.7%. The Pixel-AUROC of MDPS with $N_s = 16$ outperforms the second best DM-based method RAN by 0.6%, and shares the same values as the Pixel-AUROC of DRÆM. Although DRÆM shows the same anomaly localization performance as MDPS, the reconstruction results of DRÆM suffer from excessive blurring. Besides, DRÆM is a self-supervised method and performs badly for the real anomalies which differ significantly from the pseudo ones generated for training [Lu *et al.*, 2023]. However, MDPS does not require any pseudo training samples, thus, has higher generalization ability than DRÆM.

BTAD. To further demonstrate the generalizability and superiority of our method, we compare MDPS with several recent state-of-the-art AD methods on the BTAD dataset, including four reconstruction-based methods (DRÆM, VT-ADL [Mishra *et al.*, 2021], AnoDDPM and AutoDDPM), two representation-based methods (PatchCore [Roth *et al.*, 2022] and PaDiM [Defard *et al.*, 2021]), and a normalizing flow based method (PyramidFlow [Lei *et al.*, 2023]). The results are shown in Table 2. From Table 2, we can find that MDPS still achieves the competitive performance for anomaly detection and localization under the metric of average AUROC. Note that DRÆM shows a significant decrease performance

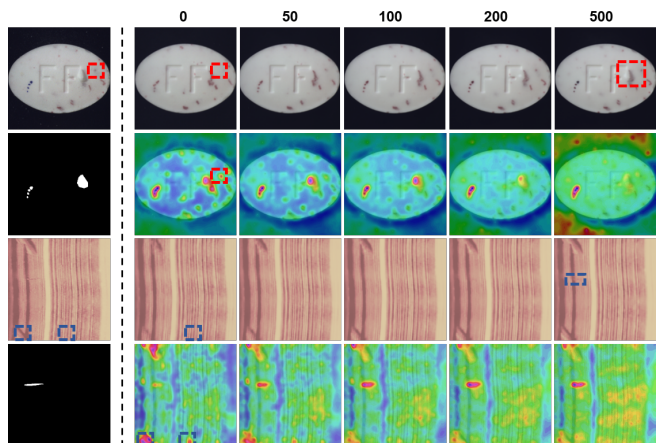


Figure 3: Qualitative comparison of different guidance scale ρ . The left side of the dotted line represents the original images and ground truths. The first and third lines on the right side of the dotted line represent the reconstructed image, and the second and fourth lines represent the heatmap. Note areas in the dotted boxes.

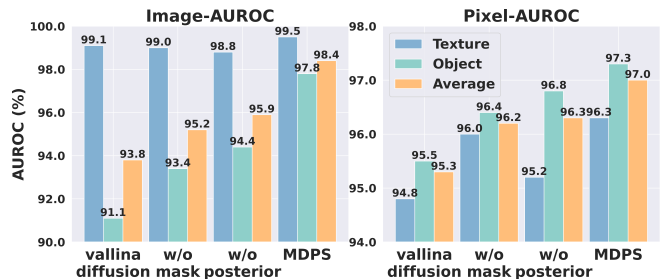


Figure 4: Ablation results of MDPS on MVTec.

on the BTAD dataset compared with the results of MVTec due to its limited generalization ability [Lu *et al.*, 2023].

Besides, we compare MDPS with several state-of-the-art AD methods on time consumption in Table 3, including PatchCore, DRÆM and AutoDDPM. Thanks to the acceleration of DDIM, MDPS shows comparable computational cost compared with AutoDDPM and PatchCore, but cannot compete with the AE-based method DRÆM.

5.2 Ablation Study

In this section, we conduct a series of ablation experiments on MVTec to discuss the selection of hyperparameters and validate the designs of MDPS and anomaly scoring.

Selection of Hyperparameters. We empirically set $T = 200$ and $N = 10$ for MDPS¹. Then, we mainly focus on selection of the guidance scale ρ for Eqn.(17), since ρ is related to anomalies and is a key hyperparameter to influence the reconstruction quality. Several normal image reconstruction examples of MVTec and BTAD are displayed to show the influence of ρ in Figure 3. In Figure 3, partial normal texture details are destructed when $\rho < 100$, which would lead to misjudgment of anomalous pixels in the normal region; however, when $\rho > 100$, anomalous texture details appear in the reconstructed image, which would lead to misjudgment

¹Selection of T and N are shown in Supplementary Material.

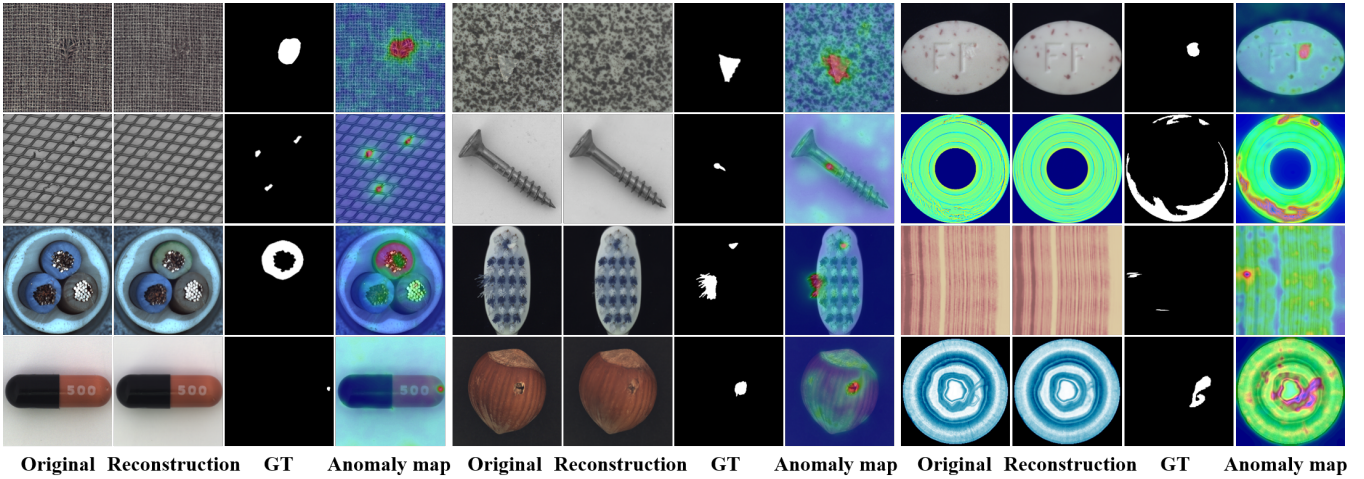


Figure 5: Qualitative results of our method. We choose 12 examples from MVTEC and BTAD and more results can be found in the supplementary. From left to right are original images, normal reconstruction images, ground truth and our localization results.

Table 4: Ablation results of the designed metric (AUROC %).

	Pixel-only	Perceptual-only	Pixel+Perceptual
Image	90.7	91.2	98.8
Pixel	90.5	92.0	97.3

Table 5: The results with different N_s on MVTEC (AUROC %).

N_s	1	2	4	8	16
Image	98.37	98.45	98.20	98.48	98.77
Pixel	96.96	97.02	97.24	97.23	97.32

of normal pixels in the anomalous region. From Figure 3, we can find that ρ can control the sensitivity for anomalies in the normal image reconstruction, and MDPS has better reconstruction quality when ρ is set to 100.

Effectiveness of MDPS. To validate the designs of MDPS, we conduct a group of ablation experiments and display the results in Figure 4. In Figure 4, “vanilla DDIM” represents only using the sampling process of DDIM, i.e., $m = 1$ and $\rho = 0$; “w/o mask” represents $m = 1$; ‘w/o posterior’ represents $\rho = 0$, i.e., the problem of normal image reconstruction is modeled as prior sampling instead of posterior sampling for normal images. From Figure 4, we can find that MDPS achieves higher values of Image-AUROC and Pixel-AUROC than MDPS with $m = 1$ or $\rho = 0$ no matter for textural or object categories, which validates the modeling of posterior sampling and the design of mask image m in MDPS.

To further validate the effectiveness of MDPS, we display qualitative results of MDPS in Figure 5. From Figure 5, we can find that MDPS can reconstruct high-quality normal images given various test images with different anomalies.

Effectiveness of anomaly scoring. The effectiveness of anomaly scoring is related to the difference metric for two images and the number of posterior samples N_s . To validate the designed metric of Eqn.(18), we compare the results of Eqn.(18) with those of only pixel-level or perceptual-level metrics in Table 4. The results in Table 4 show that the anomaly scores obtained from pixel-level and perceptual-level can improve AUROC efficiently.

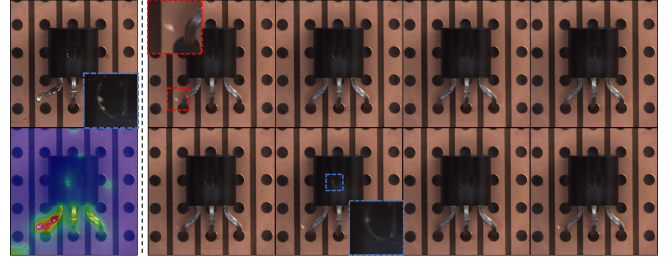


Figure 6: Multiple normal samples. **Left:** the original image and predicted heatmap. **Right:** different normal sampling results.

Then, we discuss the influence of N_s on anomaly scoring. The results with different values of N_s are shown in Table 5. It can be observed that larger value of N_s leads to further improvements of AUROC. In Figure 6, we display multiple normal samplers of a test image reconstructed by MDPS. These normal images exhibit some subtle differences. But the average of difference maps for these normal samplers can rectify misjudgment caused by a single one.

6 Conclusion

This paper proposes MDPS, a novel and highly interpretable UAD method. MDPS generates multiple normal images based on diffusion posterior sampling under Bayesian framework. Using a combination of pixel-level and perceptual-level metrics, MDPS averages all difference maps between multiple reconstructed normal images and the test image to obtain the anomaly scores accurately. Exhaustive experiments show MDPS achieves high reconstruction quality and state-of-the-art performance for anomaly detection and localization compared with recent UAD methods, including 98.8% Image-AUROC and 97.3% Pixel-AUROC on MVTEC dataset, 99.5% Image-AUROC and 97.6% Pixel-AUROC on BTAD dataset. However, MDPS suffers from high computational cost caused by diffusion posterior samplings. In future work, we would try to reduce the computational cost through knowledge distillation.

Acknowledgments

The work is supported in part by Natural Science Foundation of Guangdong Province (2022A1515010493), in part by Shenzhen Science and Technology Program (No. JCYJ20210324140407021), in part by Natural Science Foundation of China (No. 62372082) and Natural Science Foundation of Sichuan Province (No.2023NSFSC0485).

References

- [Akçay *et al.*, 2019] Samet Akçay, Amir Atapour-Abarghouei, and Toby P Breckon. Skip-GANomaly: Skip connected and adversarially trained encoder-decoder anomaly detection. In *International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2019.
- [Baur *et al.*, 2021] Christoph Baur, Stefan Denner, Benedikt Wiestler, Nassir Navab, and Shadi Albarqouni. Autoencoders for unsupervised anomaly segmentation in brain mr images: A comparative study. *Medical Image Analysis*, 69:101952, 2021.
- [Bercea *et al.*, 2023] Cosmin I. Bercea, Michael Neumayr, Daniel Rueckert, and Julia A Schnabel. Mask, stitch, and re-sample: Enhancing robustness and generalizability in anomaly detection through automatic diffusion models. In *ICML 3rd Workshop on Interpretable Machine Learning in Healthcare (IMLH)*, 2023.
- [Bergmann *et al.*, 2019] Paul Bergmann., Sindy Löwe., Michael Fauser., David Sattlegger., and Carsten Steger. Improving unsupervised defect segmentation by applying structural similarity to autoencoders. In *International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications (VISIGRAPP)*, pages 372–380. INSTICC, SciTePress, 2019.
- [Bergmann *et al.*, 2021] Paul Bergmann, Kilian Batzner, Michael Fauser, David Sattlegger, and Carsten Steger. The MVTEC anomaly detection dataset: a comprehensive real-world dataset for unsupervised anomaly detection. *International Journal of Computer Vision*, 129(4):1038–1059, 2021.
- [Chung *et al.*, 2023] Hyungjin Chung, Jeongsol Kim, Michael Thompson Mccann, Marc Louis Klasky, and Jong Chul Ye. Diffusion posterior sampling for general noisy inverse problems. In *International Conference on Learning Representations (ICLR)*, 2023.
- [Contreras-Cruz *et al.*, 2023] Marco A. Contreras-Cruz, Fernando E. Correa-Tome, Rigoberto Lopez-Padilla, and Juan-Pablo Ramirez-Paredes. Generative adversarial networks for anomaly detection in aerial images. *Computers and Electrical Engineering*, 106:108470, 2023.
- [Defard *et al.*, 2021] Thomas Defard, Aleksandr Setkov, Angélique Loesch, and Romaric Audigier. Padim: a patch distribution modeling framework for anomaly detection and localization. In *International Conference on Pattern Recognition (ICPR)*, pages 475–489. Springer, 2021.
- [Dehaene *et al.*, 2020] David Dehaene, Oriël Frigo, Sébastien Combrexelle, and Pierre Eline. Iterative energy-based projection on a normal data manifold for anomaly localization. In *International Conference on Learning Representations (ICLR)*, 2020.
- [Dhariwal and Nichol, 2021] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems (NIPS)*, 34:8780–8794, 2021.
- [Diers and Pigorsch, 2023] Jan Diers and Christian Pigorsch. A survey of methods for automated quality control based on images. *International Journal of Computer Vision*, pages 1–29, 2023.
- [Gonzalez-Jimenez *et al.*, 2023] Alvaro Gonzalez-Jimenez, Simone Lionetti, Marc Pouly, and Alexander A Navarini. Sano: Score-based diffusion model for anomaly localization in dermatology. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2987–2993, 2023.
- [Ho *et al.*, 2020] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems (NIPS)*, 33:6840–6851, 2020.
- [Lei *et al.*, 2023] Jiarui Lei, Xiaobo Hu, Yue Wang, and Dong Liu. Pyramidflow: High-resolution defect contrastive localization using pyramid normalizing flow. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14143–14152, 2023.
- [Liu *et al.*, 2020] Wenqian Liu, Runze Li, Meng Zheng, Srikrishna Karanam, Ziyang Wu, Bir Bhanu, Richard J Radke, and Octavia Camps. Towards visually explaining variational autoencoders. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8642–8651, 2020.
- [Liu *et al.*, 2023] Ruikang Liu, Weiming Liu, Zhongxing Zheng, Liang Wang, Liang Mao, Qisheng Qiu, and Guangzheng Ling. Anomaly-gan: A data augmentation method for train surface anomaly detection. *Expert Systems with Applications*, 228:120284, 2023.
- [Lu *et al.*, 2023] Fanbin Lu, Xufeng Yao, Chi-Wing Fu, and Jiaya Jia. Removing anomalies as noises for industrial defect localization. In *IEEE International Conference on Computer Vision (ICCV)*, pages 16166–16175, 2023.
- [Mishra *et al.*, 2021] Pankaj Mishra, Riccardo Verk, Daniele Fornasier, Claudio Piciarelli, and Gian Luca Foresti. Vtadl: A vision transformer network for image anomaly detection and localization. In *International Symposium on Industrial Electronics (ISIE)*, pages 01–06. IEEE, 2021.
- [Pinaya *et al.*, 2022] Walter H. L. Pinaya, Mark S. Graham, and et al. Fast unsupervised brain anomaly detection and segmentation with diffusion models. In *International Conference on Medical image computing and computer-assisted intervention*, pages 705–714. Springer, 2022.
- [Roth *et al.*, 2022] Karsten Roth, Latha Pemula, Joaquin Zepeda, Bernhard Schölkopf, Thomas Brox, and Peter Gehler. Towards total recall in industrial anomaly detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14318–14328, 2022.

- [Ruff *et al.*, 2021] Lukas Ruff, Jacob R Kauffmann, Robert A Vandermeulen, Grégoire Montavon, Wojciech Samek, Marius Kloft, Thomas G Dietterich, and Klaus-Robert Müller. A unifying review of deep and shallow anomaly detection. *Proceedings of the IEEE*, 109(5):756–795, 2021.
- [Schlegl *et al.*, 2017] 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.
- [Schlegl *et al.*, 2019] 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*, 54:30–44, 2019.
- [Song *et al.*, 2021a] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations (ICLR)*, 2021.
- [Song *et al.*, 2021b] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations (ICLR)*, 2021.
- [Song *et al.*, 2022] Jouwon Song, Kyeongbo Kong, Ye-In Park, Seong-Gyun Kim, and Suk-Ju Kang. Anomaly segmentation network using self-supervised learning. In *AAAI Workshop on AI for Design and Manufacturing (ADAM)*, 2022.
- [Sun *et al.*, 2023] Zhongju Sun, Jian Wang, and Yakun Li. Ramfae: a novel unsupervised visual anomaly detection method based on autoencoder. *International Journal of Machine Learning and Cybernetics*, 2023.
- [Tang *et al.*, 2020] Ta-Wei Tang, Wei-Han Kuo, Jauh-Hsiang Lan, Chien-Fang Ding, Hakiem Hsu, and Hong-Tsu Young. Anomaly detection neural network with dual auto-encoders gan and its industrial inspection applications. *Sensors*, 20(12):3336, 2020.
- [Wolleb *et al.*, 2022] Julia Wolleb, Florentin Bieder, Robin Sandkühler, and Philippe C Cattin. Diffusion models for medical anomaly detection. In *International Conference on Medical image computing and computer-assisted intervention*, pages 35–45. Springer, 2022.
- [Wyatt *et al.*, 2022] Julian Wyatt, Adam Leach, Sebastian M Schmon, and Chris G Willcocks. Anoddpm: Anomaly detection with denoising diffusion probabilistic models using simplex noise. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 650–656, 2022.
- [Xia *et al.*, 2022] Xuan Xia, Xizhou Pan, Nan Li, Xing He, Lin Ma, Xiaoguang Zhang, and Ning Ding. Gan-based anomaly detection: A review. *Neurocomputing*, 493:497–535, 2022.
- [You *et al.*, 2022a] Zhiyuan You, Lei Cui, Yujun Shen, Kai Yang, Xin Lu, Yu Zheng, and Xinyi Le. A unified model for multi-class anomaly detection. *Advances in Neural Information Processing Systems (NIPS)*, 35:4571–4584, 2022.
- [You *et al.*, 2022b] Zhiyuan You, Kai Yang, Wenhan Luo, Lei Cui, Yu Zheng, and Xinyi Le. ADTR: Anomaly detection transformer with feature reconstruction. In *International Conference on Neural Information Processing (ICONIP)*, pages 298–310. Springer, 2022.
- [Zagoruyko and Komodakis 2016] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *British Machine Vision Conference (BMVC)*, 2016.
- [Zavrtanik *et al.*, 2021] Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. Draem-a discriminatively trained reconstruction embedding for surface anomaly detection. In *IEEE International Conference on Computer Vision (ICCV)*, pages 8330–8339, 2021.
- [Zhang *et al.*, 2018] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 586–595, 2018.