

図2：当社のノイズ除去拡散モデルは、異常のない画像のみで訓練されます。推論時、異常サンプルに異なるスケールのノイズが追加されます。十分な大きさのノイズを加えると、異常ピクセルは正常ピクセルと区別できなくなり、再構築が容易になります。後方分布 $q(x_{t-1} | x_t, x_0)$ と推定分布 $p_{\theta}(x_{t-1} | x_t)$ の間のKLダイバージェンスをピクセルレベルの異常スコアとして採用します。特徴再構築のMSE誤差を特徴レベルのスコアとして使用します。異なるノイズスケールからの結果の平均を出力として採用します。

これにより、一連のノイズ付き画像 (x_1, x_2, \dots, x_T) が生成されます。導入されるガウスノイズの分散は $\{\beta_t\}_{t=1,2,\dots,T}$ で表されます。データ分布と追加されたノイズがともにガウス分布であるため、ノイズ付き画像 x_t の閉形式は：

$$q(x_t | x_0) = \mathcal{N}(x_t; \sqrt{\alpha_t} x_0, (1 - \alpha_t) \mathbf{I}), \quad (1)$$

ここで、 $\alpha_t = \prod_{s=1}^t (1 - \beta_s)$ です。拡散モデルは、 $p_{\theta}(x_0 : T) = \int p(x_0 : T) dx_{1:T}$ で表されます。画像生成プロセス中、モデルはまず一様ガウス分布 $p_T(x) \propto \mathcal{N}(x_T; 0, \mathbf{I})$ からサンプリングし、推定された分布 $p_{\theta}(x_t)$ からサンプリングすることで画像を徐々にノイズ除去します：

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p(x_{t-1} | x_t), \quad (2)$$

$$p_{\theta}(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_{\theta}(x_t, t)). \quad (3)$$

拡散モデルの訓練はオートエンコーダーとして扱えます。DDPM [15] で提案されたように、拡散モデルはMSE損失を用いてノイズのスケールを予測するように訓練されます。

$$L_{mse} = \mathbb{E}_{t, x_0, \epsilon} [(\epsilon - \epsilon_{\theta}(x_t, t))^2] \quad (4)$$

変分境界に基づく追加のトレーニング損失を用い、拡散モデル自体がノイズの分散を自動的に学習するように設計されています。[23]で提案された方法です：

$$L_{vbl} = L_0 + L_1 + \dots + L_{T-1} + L_T, \quad (5)$$

$$L_0 = -\log p_{\theta}(x_0 | x_1), \quad (6)$$

$$L_{t-1} = D_{KL}(q(x_{t-1} | x_t, x_0) || p_{\theta}(x_{t-1} | x_t)), \quad (7)$$

$$L_T = D_{KL}(q(x_T | x_0) || p(x_T)). \quad (8)$$

3.2. 異常検出用のノイズ除去モデル

自動エンコーダーに基づく従来の再構築手法[2, 5]は、自動エンコーダーが訓練中に同一のマッピングに退化するため、異常の再構築に成功する問題を抱えています。しかし、ノイズを含む画像での再構築はこの問題を回避します。図2に示すように、異常画像に徐々にノイズを加えると、ノイズレベルが大きな場合、異常領域が消失し、正常サンプルのピクセルと区別できなくなります。ただし、ノイズを含む画像からノイズのない画像への直接的な再構築は、重大な再構築誤差を引き起こす可能性があります。本研究では、生成型拡散モデルDDPM [15] を利用して、画像を徐々にノイズ除去し再構築します。拡散モデルは、DDPMのトレーニング手順を用いて異常のないデータでトレーニングされます。

ピクセルレベルスコア。異常検出のため、画像 x_0 にランダムなガウスノイズを加えて x_t を取得します。従来の再構築ベースの手法では、再構築画像と元の入力のRGB空間での差分を異常スコアとして使用します。しかし、このアプローチは $p(x_0 | x_t)$ の推定が困難であり、結果に大きなノイズを導入します。この制限を克服するため、事後分布 $q(x_{t-1} | x_t, x_0)$ と推定分布 $p_{\theta}(x_{t-1} | x_t)$ のKLダイバージェンスを異常スコアとして採用します。

$$s_t = KL(q(x_{t-1} | x_t, x_0) || p_{\theta}(x_{t-1} | x_t)). \quad (9)$$

図6では、KLダイバージェンスがノイズの少ない入力ピクセルの確率を適切に測定することを示しています。

特徴量レベルスコア。拡散モデルの結果は境界部で鋭いものの、頑健性がないことが観察されます

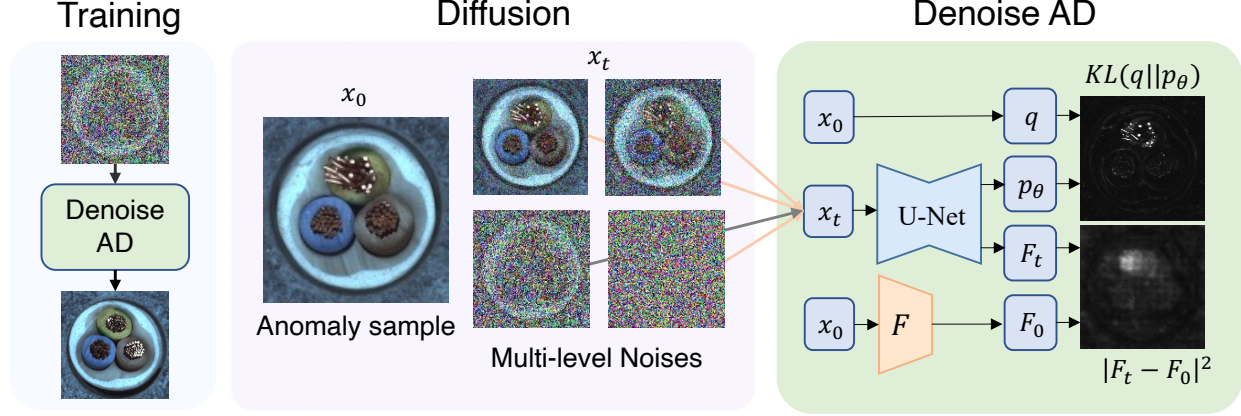


Figure 2: Our denoising diffusion model is trained with only anomaly-free images. During inference, noises of different scales are added to the anomaly sample. With large enough noises, the anomalous pixels become indistinguishable from the normal pixels and easier for reconstruction. We take the KL-divergence between the posterior distribution $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ and estimated distribution $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ as the pixel-level anomaly score. The MSE error of feature reconstruction is used as a feature-level score. We take the average of results from different noise scales as the outputs.

which leads to a series of noised images $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T)$. The variances of Gaussian noises introduced are denoted as $\{\beta_t\}_{t=1,2,\dots,T}$. Since the data distribution and noises added are both Gaussian, the closed form of a noised image \mathbf{x}_t is:

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}), \quad (1)$$

where $\bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$. The diffusion models are then represented with $p_\theta(\mathbf{x}_0) = \int p(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T}$. During the image generation process, the model first samples from uniform Gaussian distribution $p_T(\mathbf{x}) \sim \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$ and gradually denoises the image by sampling from the estimated distribution $p_\theta(\mathbf{x}_t)$:

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p(\mathbf{x}_{t-1}|\mathbf{x}_t), \quad (2)$$

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)). \quad (3)$$

The training of diffusion models can be treated as an autoencoder. As proposed in DDPM [15], the diffusion models are trained with an MSE loss to predict the scale of noises ϵ .

$$L_{mse} = \mathbb{E}_{t, \mathbf{x}_0, \epsilon} [(\epsilon - \epsilon_\theta(\mathbf{x}_t, t))^2] \quad (4)$$

An additional training loss based on the variational bound is used to automatically learn the variance of noises by the diffusion model itself, as proposed by [23]:

$$L_{vlb} = L_0 + L_1 + \dots + L_{T-1} + L_T, \quad (5)$$

$$L_0 = -\log p_\theta(\mathbf{x}_0|\mathbf{x}_1), \quad (6)$$

$$L_{t-1} = D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)||p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)), \quad (7)$$

$$L_T = D_{KL}(q(\mathbf{x}_T|\mathbf{x}_0)||p(\mathbf{x}_T)). \quad (8)$$

3.2. Denoising Model for Anomaly detection

Previous reconstruction methods based on AutoEncoder [2, 5] suffer from the successful reconstruction of anomalies because the AutoEncoder easily degrades to an identical mapping during training. However, reconstruction with noisy images prevents the issue. As illustrated in Fig 2, gradually adding noise to an anomalous image causes the anomalous regions to vanish for large noise levels, making them indistinguishable from the pixels of normal samples. Nevertheless, direct reconstruction from noisy to noise-free images can result in significant reconstruction errors. In this study, we utilize a generative diffusion model DDPM [15] to gradually denoise and reconstruct the image. The diffusion model is trained on anomaly-free data using the training procedure of DDPM.

Pixel-level score. For anomaly detection, we begin by corrupting an image \mathbf{x}_0 with random Gaussian noises to obtain \mathbf{x}_t . Previous reconstruction-based methods employ the difference between the reconstructed image and the original input in RGB space as the anomaly score. However, this approach entails a difficult estimation of $p(\mathbf{x}_0|\mathbf{x}_t)$ and introduces significant noise to the results. To address this limitation, we employ the KL divergence of the posterior distribution $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ and the estimated distribution $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ as the anomaly score,

$$s_t = KL(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)||p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)). \quad (9)$$

We show in Fig. 6 that the KL divergence correctly measures the likelihood of input pixels with much less noise.

Feature-level score. We observe that the results from the diffusion model are usually sharp in boundary but not robust