

RealNet: A Feature Selection Network with Realistic Synthetic Anomaly for Anomaly Detection

Ximiao Zhang¹ Min Xu^{1*} Xiuzhuang Zhou²

¹College of Information and Engineering, Capital Normal University

²School of Artificial Intelligence, Beijing University of Posts and Telecommunications

{2211002048, xumin}@cnu.edu.cn¹, xiuzhuang.zhou@bupt.edu.cn²

Abstract

Self-supervised feature reconstruction methods have shown promising advances in industrial image anomaly detection and localization. Despite this progress, these methods still face challenges in synthesizing realistic and diverse anomaly samples, as well as addressing the feature redundancy and pre-training bias of pre-trained feature. In this work, we introduce RealNet, a feature reconstruction network with realistic synthetic anomaly and adaptive feature selection. It is incorporated with three key innovations: First, we propose Strength-controllable Diffusion Anomaly Synthesis (SDAS), a diffusion process-based synthesis strategy capable of generating samples with varying anomaly strengths that mimic the distribution of real anomalous samples. Second, we develop Anomaly-aware Features Selection (AFS), a method for selecting representative and discriminative pre-trained feature subsets to improve anomaly detection performance while controlling computational costs. Third, we introduce Reconstruction Residuals Selection (RRS), a strategy that adaptively selects discriminative residuals for comprehensive identification of anomalous regions across multiple levels of granularity. We assess RealNet on four benchmark datasets, and our results demonstrate significant improvements in both Image AUROC and Pixel AUROC compared to the current state-of-the-art methods. The code, data, and models are available at <https://github.com/cnulab/RealNet>.

1. Introduction

Image anomaly detection is a critical task in industrial production, with wide-ranging applications in quality control and safety monitoring. While self-supervised methods [20, 32, 48, 50, 53] have gained attention for training models using synthetic anomalies, they still face challenges in synthesizing realistic and diverse anomaly images, espe-



Figure 1. SDAS generates anomaly images using only normal images. The example images are sourced from the MVTec-AD dataset [3].

cially generating complex structural anomalies and unseen anomaly categories. Due to the lack of available anomaly images and prior knowledge about anomaly categories, existing methods rely on carefully crafted data augmentation strategies [20, 32] or external data [48] for anomaly synthesis, leading to significant distribution discrepancy between synthetic anomalies and real anomalies, thereby limiting the generalization ability of anomaly detection models to real-world applications. To address these issues, we introduce Strength-controllable Diffusion Anomaly Synthesis (SDAS), a novel synthesis strategy that generates diverse samples more closely aligned with natural distributions, and offers flexibility in controlling anomaly strength. SDAS employs DDPM [16] to model the distribution of normal samples and introduces perturbation terms during the sampling process to generate samples in low probability density regions. These samples simulate various natural anomaly patterns, such as aging, structural changes, abnormal textures, and color changes, as shown in Fig. 1.

Parallel to this, feature reconstruction-based anomaly detection [8, 33, 44, 49, 53] is another promising research direction, which reconstructs the features of anomalous images as those of normal images and conducts anomaly detection and localization by reconstruction residuals. They have attracted considerable attention due to the simple

*Corresponding author.

paradigm. However, due to the high computational demands of feature reconstruction and the lack of effective feature selection strategies, existing methods either employ small-scale pre-trained CNNs [33, 44, 49] for anomaly detection or handpick layer-specific features from pre-trained network [8, 53] for reconstruction. The latest work [14] highlights the importance of feature selection, indicating that existing anomaly detection methods [30, 46] are sensitive to feature selection. The optimal pre-trained feature subset for anomaly detection varies across different categories. Therefore, devising a unified feature selection approach has become a pressing need for advancing anomaly detection. In this paper, we propose RealNet, a feature reconstruction framework that incorporates Anomaly-aware Features Selection (AFS) and Reconstruction Residuals Selection (RRS). RealNet fully exploits the discriminative capabilities of large-scale pre-trained CNNs while reducing feature redundancy and pre-training bias, enhancing anomaly detection performance while effectively controlling computational demands. For different categories, RealNet selects different pre-trained feature subsets for anomaly detection, ensuring optimal anomaly detection performance while flexibly controlling the model size. Furthermore, RealNet effectively reduces missed detections by adaptively discarding reconstruction residuals lacking anomalous information, and significantly improves the recall of anomalous regions. In summary, our contributions are fourfold:

- We propose RealNet, a feature reconstruction network that effectively leverages multi-scale pre-trained features for anomaly detection by adaptively selecting pre-trained features and reconstruction residuals. RealNet achieves state-of-the-art performance while addressing the computational cost limitations suffered by previous methods.
- We introduce Strength-controllable Diffusion Anomaly Synthesis (SDAS), a novel anomaly synthesis strategy that generates realistic and diverse anomalous samples closely aligned with natural distributions.
- We evaluate RealNet on four datasets (MVTec-AD [3], MPDD [18], BTAD [24], and VisA [55]), surpassing existing state-of-the-art methods using the same set of network architectures and hyperparameters across datasets.
- We provide the Synthetic Industrial Anomaly Dataset (SIA). SIA is generated by SDAS and consists of a total of 360,000 anomalous images from 36 categories of industrial products. SIA can be conveniently utilized for anomaly synthesis to facilitate self-supervised anomaly detection methods.

2. Related work

Unsupervised anomaly detection and localization approaches use only normal images for model training, without any anomalous data. These methods can be roughly classified into four main categories: reconstruction-based

methods [1, 2], self-supervised learning-based methods [20, 48], deep feature embedding-based methods [7, 30], and one-class classification-based methods [22, 43]. In this paper, we focus on the reconstruction-based and self-supervised learning-based methods, which are of particular relevance to our proposed RealNet framework.

Reconstruction-based methods follow a relatively consistent paradigm, which entails training a reconstruction model on normal images. The inability to effectively reconstruct anomalous regions in input images facilitates anomaly detection and localization through comparison of the original and reconstructed images. In this context, a variety of reconstruction techniques are explored, such as Autoencoder [2, 45], GAN [1, 31], Transformer [24, 28], and Diffusion model [23, 39, 52]. However, managing the reconstruction capability of the network remains challenging. In cases of complex image structures or textures, the network may produce a simplistic copy instead of selective reconstruction. Furthermore, inherent stylistic discrepancies between original and reconstructed images can lead to false positives or undetected anomalies.

Recent studies, as exemplified by [8, 33, 44, 49], have been primarily focused on anomaly detection through the reconstruction of pre-trained image features. In contrast to image-level reconstruction, multi-scale features pre-trained on ImageNet [9] demonstrate enhanced discriminative abilities to detect anomalies across a wide range of scales and diverse image patterns. However, due to the inherent feature redundancy in high-dimensional features and the pre-training bias introduced by classification tasks, the anomaly detection capability of large-scale pre-trained networks has not been fully utilized. Recent studies [33, 44, 49] use small-scale pre-trained networks to ensure controllable reconstruction costs, and other works [30, 38, 53] manually select partial layer features from pre-trained networks for anomaly detection. However, the optimal feature subset for anomaly detection varies across different categories [14], thus, these manually selecting methods often prove to be dataset-specific and suboptimal, resulting in a significant performance drop. Different from previous solutions, our RealNet presents a novel combination of efficient feature selection strategies and an optimized reconstruction process, effectively enhancing anomaly detection performance while maintaining computational efficiency.

Self-supervised learning-based methods aim to bypass the need for labels of anomalous images by setting a suitable proxy task. Notable works in this domain include CutPaste [20], which generates anomalies by transplanting image patches from one location to another, albeit with suboptimal continuity in the anomalous regions. NSA [32] uses Poisson image editing [26] for seamless image pasting to synthesize more natural anomaly regions. DRAEM [48] leverages the texture dataset DTD [5] to

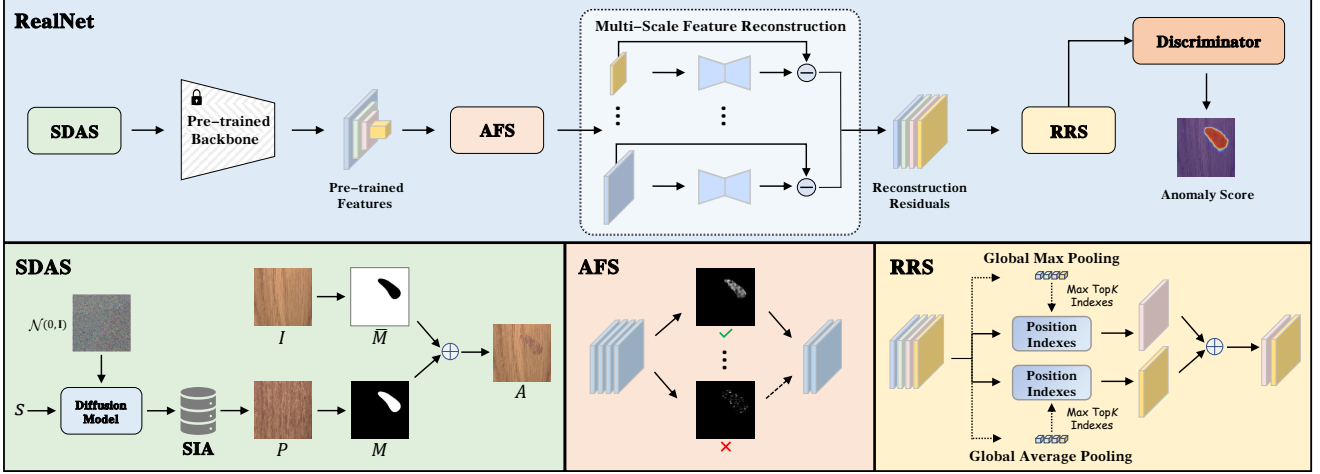


Figure 2. The pipeline of our RealNet consists of three core components: Strength-controllable Diffusion Anomaly Synthesis (SDAS), Anomaly-aware Features Selection (AFS), and Reconstruction Residuals Selection (RRS). 1) SDAS enables the synthesis of diverse, near-natural distribution anomalous images. 2) AFS refines features extracted by large-scale pre-trained CNN for dimensionality reduction. Refined features are reconstructed into corresponding normal image features by a set of reconstruction networks. 3) RRS selects reconstruction residuals most likely to identify anomalies, which are then fed into a discriminator for anomaly detection and localization.

synthesize various texture anomalies and achieve advanced self-supervised anomaly detection performance, however, it falls short when faced with specific structural anomalies, such as partial missing or misplaced elements.

The performance of self-supervised anomaly detection methods hinges on how closely the proxy task aligns with the real anomaly detection task. Anomaly synthesis, as a fundamental study in anomaly detection, has not yet received widespread exploration. Recent work [11] use StyleGAN2 [19] for image editing to generate anomalous images. However, the proposed method relies on real anomalous images and cannot generate unseen anomaly types. In contrast, SDAS operates in the probability space, free from constraints imposed by data augmentation rules or existing data, enabling effective control over anomaly strengths and the generation of realistic and diverse anomaly images using only normal images.

3. Method

In this section, we will introduce our proposed feature reconstruction framework, RealNet, which consists of three key components: Strength-controllable Diffusion Anomaly Synthesis (SDAS), Anomaly-aware Features Selection (AFS), and Reconstruction Residuals Selection (RRS). The pipeline of RealNet is illustrated in Fig. 2.

3.1. Strength-controllable Diffusion Anomaly Synthesis

Denosing Diffusion Probabilistic Models (DDPM) [16] employ a forward diffusion process to incrementally add

noise $\mathcal{N}(0, \mathbf{I})$ to the original data distribution $q(x_0)$. At time t , the conditional probability distribution of the noisy data x_t is $q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t\mathbf{I})$, where $\{\beta_t\}_{t=1}^T$ is a fixed variance schedule, and $\{x_t\}_{t=1}^T$ are the latent variables. The diffusion process is defined as a Markov chain, with joint probability distribution $q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$. Following the sum rule of Gaussian random variables, the conditional probability distribution of x_t at time t is $q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I})$, where $\alpha_t = 1 - \beta_t$, and $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$.

The reverse process is described as another Markov chain, where the mean and variance of the reverse process are parameterized by θ , i.e., $p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$. There are various ways to model $\mu_\theta(x_t, t)$; typically, neural networks $\epsilon_\theta(x_t, t)$ are used to model the noise ϵ in the diffusion process, resulting in $\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}}(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_\theta(x_t, t))$. In the training phase, our goal is to minimize the variational upper bound of the negative log-likelihood, which leads to the simplified objective:

$$\mathcal{L}_{simple} = \mathbb{E}_{t, x_0, \epsilon} [\|\epsilon - \epsilon_\theta(x_t, t)\|^2] \quad (1)$$

To generate realistic anomalous images, we first train a diffusion model to learn the distribution of normal images using Eq. (1). In reverse diffusion process characterized by $p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$, x_{t-1} is the normal image obtained at time $t - 1$. Due to the anomalous images being located in low-density regions near the normal images, we introduce an additional perturbation $s\Sigma$ to sample anomalous images, yielding $p(x'_{t-1}|x_{t-1}) = \mathcal{N}(x'_{t-1}; x_{t-1}, s\Sigma)$, where Σ is the additional introduced

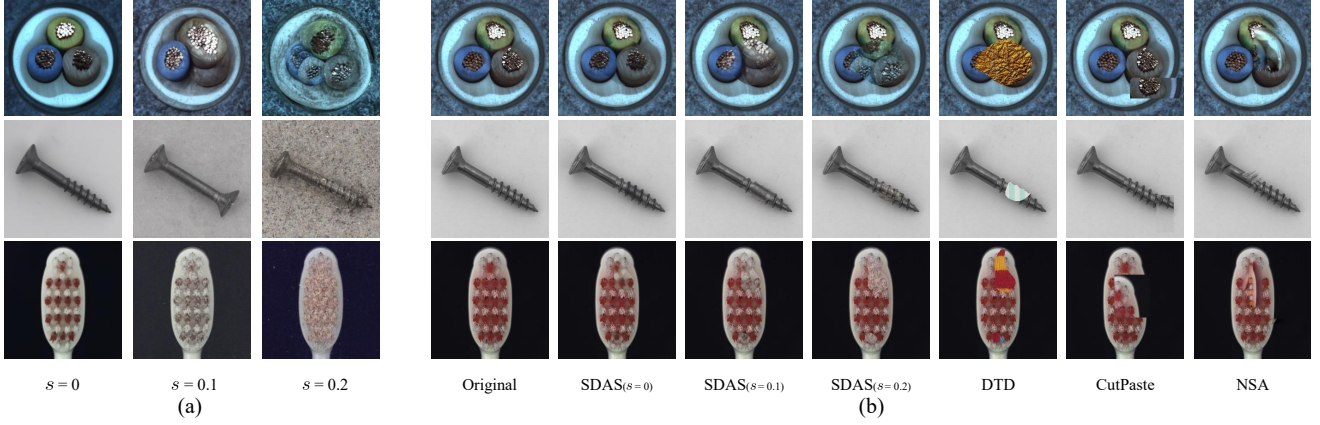


Figure 3. Anomaly image examples generated with different synthesis methods. (a) Examples generated using SDAS with different anomaly strengths s . (b) Examples featuring local anomaly regions generated by various anomaly synthesis methods.

variance, scalar s controls the anomaly strength ($s \geq 0$), and x'_{t-1} is the anomalous image obtained at time $t-1$. To simplify the anomalous synthesis process, we set $\Sigma = \Sigma_\theta(x_t, t)$, by which the conditional probability distribution of anomalous images x'_{t-1} can be written as follows:

$$p_\theta(x'_{t-1}|x_t) = \mathcal{N}(x'_{t-1}; \mu_\theta(x_t, t), (1+s)\Sigma_\theta(x_t, t)) \quad (2)$$

To ensure that the generated anomalous images are close to the distribution of normal images, we set $s \rightarrow 0$, resulting in $x'_{t-1} \approx x_{t-1}$; then we use x'_{t-1} for the next time step of the reverse diffusion process. The final form is $p_\theta(x'_{t-1}|x_t) = \mathcal{N}(x'_{t-1}; \mu_\theta(x_t, t), (1+s)\Sigma_\theta(x_t, t))$. We term this process Strength-controllable Diffusion Anomaly Synthesis (SDAS), detailed in Algorithm 1. Specifically, SDAS will generate normal images if s is set to 0.

To incorporate these anomalous images during training of anomaly detection model, we follow the approach presented in [48], utilizing a Perlin noise generator [27] to capture various anomalous shapes and binarize them into an anomaly mask M . We denote the normal image as I , the anomalous image generated by SDAS as P , and the image with local anomalies synthesized by image blending as A :

$$A = \overline{M} \odot I + (1 - \delta)(M \odot I) + \delta(M \odot P) \quad (3)$$

Algorithm 1 Strength-controllable Diffusion Anomaly Synthesis (SDAS)

Input: diffusion model $(\mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$
anomaly strength s
 $x_T \sim \mathcal{N}(0, \mathbf{I})$
for all t from T to 1 **do**
 $\mu, \Sigma \leftarrow \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)$
 $x_{t-1} \sim \mathcal{N}(\mu, (1+s)\Sigma)$
end for
return x_0

where $\overline{M} = 1 - M$, \odot denotes the element-wise multiplication operation, and δ is the opacity in the image blending. To ensure that the generated anomalous regions are located in the foreground, we use an adaptive threshold-based binarization method for foreground segmentation, similar to methods used in [32, 41, 42]. Fig. 3a shows the images generated by SDAS under different anomaly strengths, while Fig. 3b compares the images with local anomaly regions synthesized by different methods. The larger the value of s , the greater the distribution difference between the generated image and the normal image, and the more obvious the abnormal region obtained after image blending. When s is very small, imperceptible abnormal regions can be synthesized. Compared with alternative synthesis methods, the anomalies generated by SDAS are more continuous and can have very realistic structural anomalies.

3.2. Anomaly-aware Features Selection

In this section, we introduce the Anomaly-aware Features Selection (AFS) module within RealNet, a self-supervised method for pre-trained feature selection, reducing feature dimensionality and eliminating pre-training bias, as well as managing reconstruction costs. Firstly, we define a set of N triplets $\{A_n, I_n, M_n\}_{n=1}^N$, where $A_n, I_n \in R^{h \times w \times 3}$ represent anomaly images synthesized by SDAS and original normal images, and $M_n \in R^{h \times w}$ represents the corresponding anomaly mask. We denote the pre-trained network as ϕ_k , and $\phi_k(A_n) \in R^{h_k \times w_k \times c_k}$ represents the k th layer pre-trained feature extracted from A_n , where c_k represents the number of channels. For the i th feature map, $\phi_{k,i}(A_n) \in R^{h_k \times w_k}$, AFS selects m_k feature maps for reconstruction ($m_k \leq c_k$). Specifically, the feature maps indexed by k are from ResNet-like architectures, such as ResNet50 [13] or WideResNet50 [47], where $k \in \{1, 2, 3, 4\}$ represent the last layer outputs of blocks with different spatial resolutions.

For the k th layer pre-trained features, we define the following AFS loss for evaluation of the i th feature map:

$$\mathcal{L}_{AFS}(\phi_{k,i}) = \frac{1}{N} \sum_{n=1}^N \|F([\phi_{k,i}(A_n) - \phi_{k,i}(I_n)]^2) - M_n\|_2^2 \quad (4)$$

where $F(\cdot)$ is a function that performs normalization operation and aligns the resolution of $[\phi_{k,i}(A_n) - \phi_{k,i}(I_n)]^2$ to M_n . Given the feature reconstruction process for anomalous images, we train a reconstruction network to infer $\phi_{k,i}(I_n)$ based on $\phi_{k,i}(A_n)$, which enables the detection and localization of anomalies through $[\phi_{k,i}(A_n) - \phi_{k,i}(I_n)]^2$. Ideally, $[\phi_{k,i}(A_n) - \phi_{k,i}(I_n)]^2$ should closely approximate M_n . The $\mathcal{L}_{AFS}(\phi_{k,i})$ represents the capability of $\phi_{k,i}$ in identifying anomalous regions. Due to the unavailability of real anomalous samples, we employ synthetic anomalies for feature selection. For the k th layer of pre-trained features, AFS selects m_k feature maps with the smallest \mathcal{L}_{AFS} for reconstruction. We denote the AFS as $\varphi_k(\cdot)$, and $\varphi_k(A_n) \in R^{h_k \times w_k \times m_k}$, where $m_k \leq c_k$. We perform AFS on each layer of pre-trained features separately, and finally obtain selected multi-scale features $\{\varphi_1(A_n), \dots, \varphi_K(A_n)\}$. In this process, each layer's feature dimension $\{m_1, \dots, m_K\}$ serves as a set of hyperparameters. Specifically, for RealNet, AFS operation is performed only once on the pre-trained features of each layer, and the index of the selected feature maps is cached for subsequent training and inference.

AFS adaptively selects a subset of features from all available layers for anomaly detection, offering the following advantages compared to conventional methods [30, 38, 53] that select all features from partial layers: 1) AFS reduces feature redundancy within layers and mitigates pre-training bias, enhancing both feature representativeness and discriminability to improve anomaly detection performance. 2) AFS broadens the receptive field to enhance multi-scale anomaly detection capabilities. 3) AFS distinguishes the dimensions of pre-trained features from those employed for anomaly detection, ensuring efficient control over computational costs and flexible customization of the model size.

In RealNet, a set of reconstruction networks $\{G_1, \dots, G_K\}$ are designed to reconstruct the selected synthetic anomalous features $\{\varphi_1(A_n), \dots, \varphi_K(A_n)\}$ into the original image features $\{\varphi_1(I_n), \dots, \varphi_K(I_n)\}$ at various resolutions. The loss function \mathcal{L}_{recon} is defined as:

$$\mathcal{L}_{recon}(A, I) = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \|G_k(\varphi_k(A_n)) - \varphi_k(I_n)\|_2^2 \quad (5)$$

During the reconstruction process, we intentionally forgo aligning multi-scale features [30, 33, 44] to preserve optimal performance. This choice is motivated by the po-

tential drawbacks associated with aligning low-resolution features through down-sampling, which could compromise the network's detection resolution and increase the risk of misidentifying anomalies. On the other hand, aligning high-resolution features using up-sampling may result in unnecessary feature redundancy, leading to elevated reconstruction costs. A detailed discussion on the reconstruction network architectures can be found in Appendix C.

3.3. Reconstruction Residuals Selection

In this section, we present the Reconstruction Residuals Selection (RRS) module. Reconstruction residuals are denoted as $\{E_1(A_n), \dots, E_K(A_n)\}$, where $E_k(A_n) = [\varphi_k(A_n) - G_k(\varphi_k(A_n))]^2$. To obtain the global reconstruction residual $E(A_n) \in R^{h' \times w' \times m'}$, we up-sample the low-resolution reconstruction residuals and concatenate them channel-wise, where $m' = \sum_{k=1}^K m_k$, $h' = \max(h_1, \dots, h_K)$, and $w' = \max(w_1, \dots, w_K)$.

The reconstruction residuals in $E(A_n)$ is obtained from the pre-trained features of reconstructing corresponding layer, and the features of the same resolution only have good ability to capture anomalies within a certain range. For instance, subtle low-level texture anomalies can be effectively captured exclusively by reconstruction residuals derived from low-level features. Therefore, RRS selects only a subset of reconstruction residuals that contain the most anomalous information for the anomaly score generation, to achieve the highest possible recall of anomalous regions.

Firstly, RRS performs GlobalMaxPooling (GMP) and GlobalAveragePooling (GAP) on $E(A_n)$ to obtain $E_{GMP}(A_n), E_{GAP}(A_n) \in R^{m'}$ respectively. The r largest elements in $E_{GMP}(A_n)$ and $E_{GAP}(A_n)$ are then used to index the positions of $E(A_n)$ and obtain $E_{max}(A_n, r), E_{avg}(A_n, r) \in R^{h' \times w' \times r}$, which respectively represent the TopK reconstruction residuals with the highest maximum and average values. To avoid missed detections caused by inadequate resolution, reconstruction residuals with insufficient anomalous information are discarded in RRS.

As GMP and GAP respectively represent local and global properties spatially, E_{max} is more effective in capturing local anomalies in small areas, while E_{avg} focuses on selecting anomalies with large spans. Combining E_{max} and E_{avg} together can enhance the RRS's ability to capture anomalies of various scales. We define the RRS operator as $E_{RRS}(A_n, r) \in R^{h' \times w' \times r}$. $E_{RRS}(A_n, r)$ concatenates $E_{max}(A_n, r/2)$ and $E_{avg}(A_n, r/2)$. Finally, we feed the $E_{RRS}(A_n, r)$ into a discriminator, which maps the reconstruction residual to the image-level resolution, obtaining the final anomaly scores. The maximum value in anomaly scores is used as the image-level anomaly score. We use cross entropy loss $\mathcal{L}_{seg}(A, M)$ to supervise the training of

Table 1. Comparison of SIA with alternative anomaly synthesis approaches on the MVTec-AD dataset [3], employing Image AUROC (%), Pixel AUROC (%), and PRO (%) as evaluation metrics.

Category		SIA	DTD [5]	NSA [32]	CutPaste [20]
Texture	Carpet	(99.84, 99.19, 96.41)	(100.0, 99.27, 96.96)	(99.80, 98.60, 88.77)	(99.24, 98.42, 93.85)
	Grid	(100.0, 99.51, 97.28)	(100.0, 99.57, 97.14)	(100.0, 99.32, 91.31)	(100.0, 99.18, 92.53)
	Leather	(100.0, 99.76, 96.22)	(100.0, 99.77, 96.41)	(100.0, 99.24, 96.85)	(100.0, 99.41, 92.12)
	Tile	(99.96, 99.44, 97.70)	(100.0, 99.35, 95.27)	(100.0, 97.40, 86.45)	(99.86, 97.63, 84.39)
	Wood	(99.21, 98.22, 90.54)	(99.65, 98.28, 91.23)	(97.63, 93.30, 87.20)	(98.95, 95.29, 81.47)
	AVG	(99.80, 99.22, 95.63)	(99.93, 99.25, 95.40)	(99.49, 97.57, 90.11)	(99.61, 97.99, 88.87)
Object	Bottle	(100.0, 99.30, 95.62)	(100.0, 99.35, 95.57)	(100.0, 99.37, 93.49)	(100.0, 99.14, 91.41)
	Cable	(99.19, 98.10, 93.38)	(98.95, 97.84, 90.36)	(99.33, 97.62, 93.26)	(96.35, 96.23, 86.05)
	Capsule	(99.56, 99.32, 84.48)	(99.32, 99.19, 82.28)	(99.04, 99.27, 85.77)	(98.48, 99.10, 79.55)
	hazelnut	(100.0, 99.68, 93.14)	(100.0, 99.46, 93.46)	(100.0, 99.25, 94.41)	(100.0, 99.03, 91.51)
	Metal Nut	(99.76, 98.58, 94.39)	(99.90, 98.58, 96.49)	(100.0, 99.11, 93.27)	(99.90, 98.03, 89.69)
	Pill	(99.13, 99.02, 91.04)	(98.36, 98.88, 84.44)	(97.19, 98.28, 95.15)	(97.22, 98.96, 86.48)
	Screw	(98.83, 99.45, 87.90)	(97.72, 99.36, 85.22)	(98.79, 99.62, 93.74)	(92.74, 98.53, 79.63)
	Toothbrush	(99.44, 98.71, 91.57)	(99.44, 98.69, 90.87)	(100.0, 99.18, 89.20)	(99.17, 98.85, 78.48)
	Transistor	(100.0, 98.00, 92.92)	(99.71, 97.15, 86.56)	(98.54, 95.67, 79.09)	(99.38, 96.32, 76.52)
	zipper	(99.82, 99.17, 93.43)	(99.68, 99.02, 88.77)	(99.90, 98.91, 93.05)	(99.61, 98.03, 92.26)
	AVG	(99.57, 98.93, 91.79)	(99.31, 98.75, 89.40)	(99.28, 98.63, 91.04)	(98.29, 98.22, 85.16)
AVG		(99.65, 99.03, 93.07)	(99.52, 98.92, 91.40)	(99.35, 98.28, 90.73)	(98.73, 98.14, 86.40)

Table 2. Comparison of RealNet with alternative anomaly detection methods on the MVTec-AD dataset [3].

Metric	PatchCore [30]	SimpleNet [21]	FastFlow [46]	DRAEM+SSPCAB [29]	DSR [49]	UniAD [44]	RD++ [38]	DeSTSeg [53]	DiffAD [52]	RealNet
Image AUROC	99.1	99.6	99.3	98.9	98.2	96.6	99.4	98.6	98.7	99.6
Pixel AUROC	98.1	98.1	98.1	97.2	-	96.6	98.3	97.9	98.3	99.0

discriminator. The overall loss of RealNet is:

$$\mathcal{L}(A, I, M) = \mathcal{L}_{recon}(A, I) + \mathcal{L}_{seg}(A, M) \quad (6)$$

3.4. Synthetic Industrial Anomaly Dataset

To facilitate the reuse of generated anomaly images by SDAS, we constructed the Synthetic Industrial Anomaly Dataset (SIA). SIA comprises anomaly images for 36 categories from four industrial anomaly detection datasets, including MVTec-AD [3], MPDD [18], BTAD [24], and VisA [55]. We generated 10,000 anomaly images with a resolution of 256×256 for each category, with anomaly strength s uniformly sampled between 0.1 and 0.2. SIA can be conveniently used for synthesizing anomaly images through image blending, as described in Eq. (3), and can serve as an effective alternative to the widely used DTD dataset [5].

4. Experiment

4.1. Experimental setup

Datasets. We conduct extensive evaluations on four datasets, including MVTec-AD [3], MPDD [18], BTAD [24], and VisA [55]. MVTec-AD [3] contains 5,354 images from 15 categories for industrial anomaly detection tasks, including 10 object categories and 5 texture categories. MPDD [18] contains 1,346 images from 6 types of industrial metal products with varying lighting conditions, non-uniform backgrounds, and multiple products in

each image. Furthermore, the placement orientation, shooting distance, and position of the products are also varied. BTAD [24] contains images of 3 industrial products from the real world. VisA [55] is comprised of 9,621 normal images and 1,200 anomaly images from 12 categories. Certain categories demonstrate intricate structures, as exemplified by PCBs, while others consist of multiple objects that require detection, such as Capsules, thus rendering detection and localization a challenging task.

Metrics. To evaluate the performance of image-level anomaly detection, we use the Area Under the Receiver Operator Curve (AUROC) metric, as in previous works [3, 18, 24, 55]. For pixel-level anomaly location, we use Pixel AUROC and Per Region Overlap (PRO) [4].

Implementation details. We evaluate RealNet on four datasets with consistent network architectures and hyperparameters, without specific tuning for individual categories. We use a WideResNet50 [47] pre-trained on ImageNet [9] as the backbone. In AFS, we set the dimension of pre-trained feature of each layer to $\{256, 512, 512, 256\}$ for reconstruction. For RRS, 1/3 of the reconstruction residuals are reserved to generate the final anomaly scores. For SDAS, we train the diffusion model following [10] and use the SIA dataset for anomaly synthesis. Both SDAS and anomaly detection are performed at a resolution of 256×256 without center cropping, with a batch size of 16, and we use 64 batches of synthetic anomaly images for AFS. More details can be found in Appendix B.

Table 3. Comparison of SIA with DTD [5] and CutPaste [20] on the MPDD dataset [18], employing Image AUROC (%), Pixel AUROC (%), and PRO (%) as evaluation metrics.

Category	SIA	DTD [5]	CutPaste [20]
Bracket Black	(94.95, 99.27, 87.10)	(89.49, 98.90, 88.57)	(66.42, 96.67, 56.53)
Bracket Brown	(96.83, 97.81, 94.36)	(92.99, 97.35, 92.64)	(95.48, 97.54, 55.17)
Bracket White	(88.78, 97.44, 84.00)	(86.67, 98.59, 77.08)	(88.44, 96.51, 64.32)
Connector	(100.0, 97.46, 84.79)	(99.05, 97.76, 65.91)	(99.05, 98.47, 74.05)
Metal Plate	(100.0, 99.28, 94.44)	(100.0, 99.35, 93.78)	(99.95, 98.83, 92.69)
Tubes	(97.51, 97.94, 93.29)	(92.62, 99.01, 96.49)	(91.49, 98.09, 92.99)
AVG	(96.35, 98.20, 89.66)	(93.47, 98.49, 85.75)	(90.14, 97.69, 72.63)

Table 4. Comparison of RealNet with alternative anomaly detection methods on the MPDD dataset [18].

Metric	PatchCore [30]	CFlow [12]	PaDiM [7]	SPADE [6]	DAGAN [36]	Skip-GANomaly [1]	RealNet
Image AUROC	82.1	86.1	74.8	77.1	72.5	64.8	96.3
Pixel AUROC	95.7	97.7	96.7	95.9	83.3	82.2	98.2

4.2. Anomaly detection on MVTec-AD

We train RealNet using SIA and alternative anomaly synthetic methods on the MVTec-AD dataset [3], to evaluate the model’s performance in anomaly detection and localization. These methods include: 1) **DTD** [5]: This method utilizes the DTD dataset [5] to blend images with generated anomalous textures, and the data augmentation strategy in [48] is employed during the blending process. 2) **NSA** [32]: This method employs Poisson image editing [26] to seamless image editing, following parameter setting in [32]. 3) **CutPaste** [20]: This method involves random image cropping and pasting to synthesize anomaly regions.

The experimental results are shown in Tab. 1. SDAS demonstrates flexibility in controlling the anomaly strength and generates synthetic anomalies with multiple anomaly patterns, especially for the object categories, where it achieves the best detection and localization performance. Compared to other methods, SDAS is not constrained by data augmentation rules or external data, enabling the synthesis of more natural and rich functional anomalies, as shown in Fig. 3. RealNet trained using SIA achieves remarkable performance on the MVTec-AD dataset [3], with an Image AUROC of 99.65%, a Pixel AUROC of 99.03%, and a PRO score of 93.07%. Fig. 4 presents the qualitative anomaly localization results of RealNet on the MVTec-AD dataset [3]. The method exhibits remarkable pixel-level anomaly localization, proficiently identifying diverse anomaly patterns at various scales. Furthermore, RealNet can achieve a rapid inference speed of 31.93 FPS when using a single Nvidia GeForce RTX 3090, and it can perform inference using only 4GB of GPU memory. A detailed computational efficiency analysis can be found in Appendix C.

We also compare RealNet with several state-of-the-art anomaly detection methods, and the results are shown in Tab. 2. Built on the same pre-trained network, RealNet outperforms the state-of-the-art alternatives, including Deep

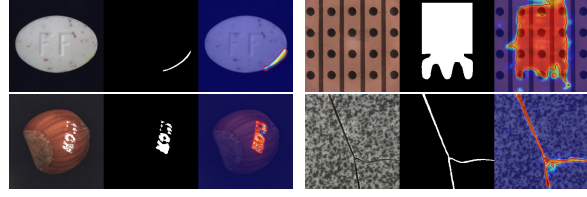


Figure 4. Qualitative results of RealNet on the MVTec-AD dataset [3]. Within each group, from left to right, are the anomaly image, ground-truth, and predicted anomaly score.

Feature Embedding-Based method (PatchCore [30] and SimpleNet [21]) and the NF-Based method (FastFlow [46]). When compared to previous reconstruction-based methods, RealNet achieves significant performance improvement.

4.3. Anomaly detection on MPDD

We evaluate RealNet on MPDD dataset [18] with SIA, DTD [5], and CutPaste [20], and the results are shown in Tab. 3. Notably, RealNet trained with SIA achieves a significant improvement of 2.88% in Image AUROC over DTD [5]. Tab. 4 shows the Image AUROC and the Pixel AUROC of RealNet and other methods on the MPDD dataset [18]. RealNet achieves an Image AUROC of 96.3%, surpassing the current best performance (CFlow [12]) by 10.2%, even without any dataset-specific tuning.

4.4. Anomaly detection on other benchmarks

To comprehensively assess the effectiveness of RealNet, we conduct experiments on the BTAD [24] and VisA [55] datasets. On the VisA dataset [55], characterized by complex structures and multiple detection objects, RealNet employing SIA yields a significant performance improvement, achieving an Image AUROC of 97.8% and a Pixel AUROC of 98.8%. In the case of the BTAD dataset [24], RealNet with SIA achieves competitive results, securing an Image AUROC of 96.1% and a Pixel AUROC of 97.9%. Detailed results can be found in Appendix C.

4.5. Ablation studies

To evaluate the effectiveness of each module of RealNet, we conduct comprehensive ablation studies on MVTec-AD dataset [3]. First, we evaluate the impact of AFS and RRS on the performance of RealNet.

W/O AFS: We replace AFS with two alternative dimensionality reduction methods, namely Random Dimensionality Reduction (RDR) [7] and Random Linear Projections

Table 5. Ablation studies of RealNet on the MVTec-AD dataset [3].

(a) The impact of AFS and RRS on RealNet.						(b) The impact of anomaly strengths on RealNet.				
	AFS	RRS	Image AUROC	Pixel AUROC	PRO	Metric	$s = 0$	$s = 0.1$	$s = 0.2$	$s \in [0.1, 0.2]$
1	-	-	94.46 / 95.67	93.38 / 95.84	79.81 / 82.26	Image AUROC	99.35	99.65	99.61	99.65
2	✓	-	96.86	96.32	84.13	Pixel AUROC	98.85	98.96	98.95	99.03
3	-	✓	99.39 / 99.09	98.66 / 98.22	92.01 / 88.38	PRO	91.80	92.16	89.36	93.07
4	✓	✓	99.65	99.03	93.07					

Figure 5. Performance of RealNet on MVTec-AD dataset [3] under various modes of reconstruction residuals selection (Max, Avg, and Max&Avg) and varying retention ratio P of reconstruction residuals.

Reduction (RLPR) [30, 40]. RDR randomly selects some dimensional features from high-dimensional features, while RLPR employs an untrained linear transformation layer for linear projection. We report the results of our RealNet with RDR and RLPR, respectively, as shown in the experiments 1 and 3 in Tab. 5a. **W/O RRS:** We feed the global reconstruction residual $E(A_n)$ into the discriminator to generate anomaly scores, and the results are shown in the experiments 1 and 2 in Tab. 5a.

As indicated by the ablation results in Tab. 5a, RRS contributes significantly to performance improvement. Using all the reconstruction residuals to generate anomaly scores, reconstruction residuals lacking of anomaly information can lead to missed anomaly regions, resulting in a significant decrease in anomaly detection performance. Furthermore, AFS yields better anomaly detection results compared to RDR and RLPR. A straightforward visualization result about AFS is provided in Appendix D.

We further investigate the impact of anomaly strength s in SDAS, and the results are shown in Tab. 5b. When s equals 0, SDAS generates normal images in high probability density regions. Blending images may introduce false positive anomaly regions, which lowers the reconstruction difficulty and confuses the discriminator, leading to sub-optimal performance. Conversely, when s is too large, the synthetic anomaly images deviate from the true distribution of anomalous images, causing RealNet’s performance to deteriorate. Our findings indicate that uniformly sampling s within a specific range serves as a robust approach for generating anomalous images. This method enables RealNet to cover a wider range of anomalous patterns, ultimately improving the overall anomaly detection performance.

In Fig. 5, we report the impact of different RRS modes

and retention ratios on the performance of RealNet. Since the setting of r is related to m_k , we introduce the retention ratio P , defined as: $P = \frac{r}{\sum_{k=1}^K m_k}$. Compared to Max and Avg modes, Max&Avg mode demonstrates superior robustness in detecting anomalies across various scales. At equal retention rates, the Max&Avg mode discards more reconstruction residuals lacking anomalous information than the Max and Avg modes, mitigating performance degradation and further emphasizing the effectiveness of the Max&Avg mode in enhancing RealNet’s anomaly detection capabilities. More ablation experiments and analysis can be found in Appendix C.

5. Conclusion

In this work, we introduce RealNet, an innovative self-supervised anomaly detection framework. Our approach integrates three core components: Strength-controllable Diffusion Anomaly Synthesis (SDAS), Anomaly-aware Features Selection (AFS), and Reconstruction Residuals Selection (RRS). These components synergistically contribute to RealNet, enabling effective exploitation of large-scale pre-trained models in anomaly detection while keeping the computational overhead within a reasonably low and acceptable range. RealNet provides a flexible foundation for future research in anomaly detection utilizing pre-trained feature reconstruction techniques. Through extensive experiments, we illustrate RealNet’s proficiency in addressing diverse real-world anomaly detection challenges.

Acknowledgements. This work was supported in part by the National Natural Science Foundation of China under Grant 62177034 and Grant 61972046.

References

- [1] Samet Akçay, Amir Atapour-Abarghouei, and Toby P Breckon. Skip-ganomaly: Skip connected and adversarially trained encoder-decoder anomaly detection. In *2019 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2019. 2, 7
- [2] Christoph Baur, Benedikt Wiestler, Shadi Albarqouni, and Nassir Navab. Deep autoencoding models for unsupervised anomaly segmentation in brain mr images. In *Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries: 4th International Workshop, BrainLes 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16, 2018, Revised Selected Papers, Part I 4*, pages 161–169. Springer, 2019. 2
- [3] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. Mvtec-ad: A comprehensive real-world dataset for unsupervised anomaly detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9592–9600, 2019. 1, 2, 6, 7, 8, 13, 14, 15, 16, 17, 18, 19, 20
- [4] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. Uninformed students: Student-teacher anomaly detection with discriminative latent embeddings. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4183–4192, 2020. 6
- [5] Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3606–3613, 2014. 2, 6, 7, 12, 13, 15, 16
- [6] Niv Cohen and Yedid Hoshen. Sub-image anomaly detection with deep pyramid correspondences. *arXiv preprint arXiv:2005.02357*, 2020. 7, 13
- [7] Thomas Defard, Aleksandr Setkov, Angelique Loesch, and Romaric Audigier. Padim: a patch distribution modeling framework for anomaly detection and localization. In *Pattern Recognition. ICPR International Workshops and Challenges: Virtual Event, January 10–15, 2021, Proceedings, Part IV*, pages 475–489. Springer, 2021. 2, 7, 15
- [8] Hanqiu Deng and Xingyu Li. Anomaly detection via reverse distillation from one-class embedding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9737–9746, 2022. 1, 2
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 2, 6, 12, 15
- [10] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. 6
- [11] Yuxuan Duan, Yan Hong, Li Niu, and Liqing Zhang. Few-shot defect image generation via defect-aware feature manipulation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 571–578, 2023. 3, 16
- [12] Denis Gudovskiy, Shun Ishizaka, and Kazuki Kozuka. Cflow-ad: Real-time unsupervised anomaly detection with localization via conditional normalizing flows. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 98–107, 2022. 7
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4, 13, 14, 15
- [14] Lars Heckler, Rebecca König, and Paul Bergmann. Exploring the importance of pretrained feature extractors for unsupervised anomaly detection and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2916–2925, 2023. 2
- [15] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 15
- [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 1, 3
- [17] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. pmlr, 2015. 12
- [18] Stepan Jezek, Martin Jonak, Radim Burget, Pavel Dvorak, and Milos Skotak. Deep learning-based defect detection of metal parts: evaluating current methods in complex conditions. In *2021 13th International Congress on Ultra Modern Telecommunications and Control Systems and Workshops (ICUMT)*, pages 66–71. IEEE, 2021. 2, 6, 7, 17, 18, 19, 20
- [19] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. 3
- [20] Chun-Liang Li, Kihyuk Sohn, Jinsung Yoon, and Tomas Pfister. Cutpaste: Self-supervised learning for anomaly detection and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9664–9674, 2021. 1, 2, 6, 7, 15, 16
- [21] Zhikang Liu, Yiming Zhou, Yuansheng Xu, and Zilei Wang. Simplenet: A simple network for image anomaly detection and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20402–20411, 2023. 6, 7
- [22] Philipp Liznerski, Lukas Ruff, Robert A. Vandermeulen, Billy Joe Franks, Marius Kloft, and Klaus Robert Muller. Explainable deep one-class classification. In *International Conference on Learning Representations*, 2021. 2
- [23] Fanbin Lu, Xufeng Yao, Chi-Wing Fu, and Jiaya Jia. Removing anomalies as noises for industrial defect localization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 16166–16175, 2023. 2
- [24] Pankaj Mishra, Riccardo Verk, Daniele Fornasier, Claudio Picciarelli, and Gian Luca Foresti. Vt-adl: A vision transformer network for image anomaly detection and localization

- tion. In *2021 IEEE 30th International Symposium on Industrial Electronics (ISIE)*, pages 01–06. IEEE, 2021. 2, 6, 7, 12, 13, 17, 18, 19, 20
- [25] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. 12
- [26] Patrick Pérez, Michel Gangnet, and Andrew Blake. Poisson image editing. In *ACM SIGGRAPH 2003 Papers*, pages 313–318. 2003. 2, 7
- [27] Ken Perlin. An image synthesizer. *ACM Siggraph Computer Graphics*, 19(3):287–296, 1985. 4
- [28] Jonathan Pirnay and Keng Chai. Inpainting transformer for anomaly detection. In *Image Analysis and Processing–ICIAP 2022: 21st International Conference, Lecce, Italy, May 23–27, 2022, Proceedings, Part II*, pages 394–406. Springer, 2022. 2
- [29] Nicolae-Cătălin Ristea, Neelu Madan, Radu Tudor Ionescu, Kamal Nasrollahi, Fahad Shabbaz Khan, Thomas B Moeslund, and Mubarak Shah. Self-supervised predictive convolutional attentive block for anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13576–13586, 2022. 6
- [30] Karsten Roth, Latha Pemula, Joaquin Zepeda, Bernhard Schölkopf, Thomas Brox, and Peter Gehler. Towards total recall in industrial anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14318–14328, 2022. 2, 5, 6, 7, 8, 13
- [31] 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 *Information Processing in Medical Imaging: 25th International Conference, IPMI 2017, Boone, NC, USA, June 25-30, 2017, Proceedings*, pages 146–157. Springer, 2017. 2
- [32] Hannah M Schlüter, Jeremy Tan, Benjamin Hou, and Bernhard Kainz. Natural synthetic anomalies for self-supervised anomaly detection and localization. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXI*, pages 474–489. Springer, 2022. 1, 2, 4, 6, 7, 15, 16
- [33] Yong Shi, Jie Yang, and Zhiquan Qi. Unsupervised anomaly segmentation via deep feature reconstruction. *Neurocomputing*, 424:9–22, 2021. 1, 2, 5, 14
- [34] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021. 12, 13
- [35] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019. 13, 14
- [36] 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. 7
- [37] Xian Tao, Dapeng Zhang, Wenzhi Ma, Zhanxin Hou, Zhen-Feng Lu, and Chandranath Adak. Unsupervised anomaly detection for surface defects with dual-siamese network. *IEEE Transactions on Industrial Informatics*, 18(11):7707–7717, 2022. 14
- [38] Tran Dinh Tien, Anh Tuan Nguyen, Nguyen Hoang Tran, Ta Duc Huy, Soan Duong, Chanh D Tr Nguyen, and Steven QH Truong. Revisiting reverse distillation for anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24511–24520, 2023. 2, 5, 6, 13
- [39] Julian Wyatt, Adam Leach, Sebastian M Schmon, and Chris G Willcocks. Anoddpm: Anomaly detection with denoising diffusion probabilistic models using simplex noise. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 650–656, 2022. 2
- [40] Jiang Xi, Jianlin Liu, Jinbao Wang, Qiang Nie, WU Kai, Yong Liu, Chengjie Wang, and Feng Zheng. Softpatch: Unsupervised anomaly detection with noisy data. In *Advances in Neural Information Processing Systems*. 8
- [41] Minghui Yang, Peng Wu, and Hui Feng. Memseg: A semi-supervised method for image surface defect detection using differences and commonalities. *Engineering Applications of Artificial Intelligence*, 119:105835, 2023. 4
- [42] Xincheng Yao, Ruqi Li, Jing Zhang, Jun Sun, and Chongyang Zhang. Explicit boundary guided semi-push-pull contrastive learning for supervised anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24490–24499, 2023. 4
- [43] Jihun Yi and Sungroh Yoon. Patch svdd: Patch-level svdd for anomaly detection and segmentation. In *Proceedings of the Asian Conference on Computer Vision*, 2020. 2, 13
- [44] Zhiyuan You, Lei Cui, Yujun Shen, Kai Yang, Xin Lu, Yu Zheng, and Xinyi Le. A unified model for multi-class anomaly detection. In *Advances in Neural Information Processing Systems*, 2022. 1, 2, 5, 6, 14, 15
- [45] Sanyapong Youkachen, Miti Ruchanurucks, Teera Phatpomnant, and Hirohiko Kaneko. Defect segmentation of hot-rolled steel strip surface by using convolutional auto-encoder and conventional image processing. In *2019 10th International Conference of Information and Communication Technology for Embedded Systems (IC-ICTES)*, pages 1–5. IEEE, 2019. 2
- [46] Jiawei Yu, Ye Zheng, Xiang Wang, Wei Li, Yushuang Wu, Rui Zhao, and Liwei Wu. Fastflow: Unsupervised anomaly detection and localization via 2d normalizing flows. *arXiv preprint arXiv:2111.07677*, 2021. 2, 6, 7, 13
- [47] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *Proceedings of the British Machine Vision Conference (BMVC)*, pages 87.1–87.12. BMVA Press, 2016. 4, 6, 13, 14, 21
- [48] Vitjan Zavrtanik, Matej Kristan, and Danijel Škočaj. Draem: A discriminatively trained reconstruction embedding for surface anomaly detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8330–8339, 2021. 1, 2, 4, 7, 13, 15, 16

- [49] Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. Dsr: A dual subspace re-projection network for surface anomaly detection. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXI*, pages 539–554. Springer, 2022. 1, 2, 6
- [50] Hui Zhang, Zuxuan Wu, Zheng Wang, Zhineng Chen, and Yu-Gang Jiang. Prototypical residual networks for anomaly detection and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16281–16291, 2023. 1
- [51] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 15, 16
- [52] Xinyi Zhang, Naiqi Li, Jiawei Li, Tao Dai, Yong Jiang, and Shu-Tao Xia. Unsupervised surface anomaly detection with diffusion probabilistic model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6782–6791, 2023. 2, 6
- [53] Xuan Zhang, Shiyu Li, Xi Li, Ping Huang, Jiulong Shan, and Ting Chen. Destseg: Segmentation guided denoising student-teacher for anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3914–3923, 2023. 1, 2, 5, 6
- [54] Ying Zhao. Omnia: A unified cnn framework for unsupervised anomaly localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3924–3933, 2023. 15
- [55] Yang Zou, Jongheon Jeong, Latha Pemula, Dongqing Zhang, and Onkar Dabeer. Spot-the-difference self-supervised pre-training for anomaly detection and segmentation. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXX*, pages 392–408. Springer, 2022. 2, 6, 7, 12, 13, 17, 18, 19, 20

RealNet: A Feature Selection Network with Realistic Synthetic Anomaly for Anomaly Detection

Supplementary Material

A. Overview

We organize this supplementary material into the following sections: Appendix B provides additional implementation details for RealNet. Appendix C provides detailed results on the BTAD [24] and VisA [55] datasets, supplementary ablation study results, an analysis of RealNet’s computational efficiency, anomaly detection results in multi-class setting, as well as synthetic anomaly image quality assessment results. Appendix D offers additional visualization results, including qualitative results of RealNet in anomaly localization, images generated by SDAS, and a straightforward visualization result of AFS. Appendix E discusses the limitations of our method.

B. More details

In SDAS, we use the learnable reverse diffusion variance [25] as $\Sigma_\theta(x_t, t)$, given by:

$$\Sigma_\theta(x_t, t) = \exp(v \log \beta_t + (1 - v) \log \tilde{\beta}_t) \quad (S1)$$

Here, β_t represents the variance of the diffusion process, while $\tilde{\beta}_t$ represents the variance of the conditional posterior distribution $q(x_{t-1}|x_t, x_0)$, and $\tilde{\beta}_t = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t} \beta_t$. The vector v is predicted by the model and weighted with β_t and $\tilde{\beta}_t$ in the \log space. We optimize $\mu_\theta(x_t, t)$ and $\Sigma_\theta(x_t, t)$ with the loss \mathcal{L}_{hybrid} :

$$\mathcal{L}_{hybrid} = \mathcal{L}_{simple} + \gamma \mathcal{L}_{vlb} \quad (S2)$$

where

$$\begin{aligned} \mathcal{L}_{vlb} &= \mathcal{L}_0 + \mathcal{L}_1 + \dots + \mathcal{L}_{T-1} + \mathcal{L}_T \\ \mathcal{L}_0 &= -\log p_\theta(x_0|x_1) \\ \mathcal{L}_{t-1} &= D_{KL}(q(x_{t-1}|x_t, x_0)||p_\theta(x_{t-1}|x_t)) \\ \mathcal{L}_T &= D_{KL}(q(x_T|x_0)||p(x_T)) \end{aligned} \quad (S3)$$

We set γ to 0.001 in Eq. (S2), and stop the gradient of $\mu_\theta(x_t, t)$ in \mathcal{L}_{vlb} during the training phase. To accelerate the convergence of the diffusion model, we initialize it with weights pre-trained on ImageNet [9]. We set the reverse diffusion step T of 20, and generating 10,000 images at a resolution of 256×256 takes 6 hours using a single NVIDIA GeForce RTX 3090.

The SDAS with DDIM [34] is described in Algorithm S1, which provides three options for applying perturbation variance in the deterministic reverse diffusion process: $\Sigma = \beta_t$, $\Sigma = \tilde{\beta}_t$, and $\Sigma = \Sigma_\theta(x_t, t)$. Experimental

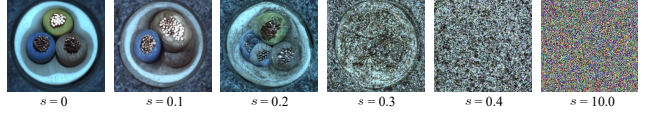


Figure S1. Sample anomaly images generated by SDAS with different anomaly strengths s .

observations show that the anomaly images obtained by ID-DPM [25] are slightly better than those obtained by DDIM [34], and therefore, we use IDDPM [25] for SDAS. Some examples can be found in Appendix D.

Fig. S1 presents examples of images generated by SDAS with a broader range of anomaly strengths. As the anomaly strength increases, the generated anomalous images contain more noise, reducing their authenticity. In the experiments, we set the anomaly strength between 0.1 and 0.2, allowing SDAS to encompass a wider range of real-world anomalies.

In RRS, the global reconstruction residual $E(A_n)$ originates from distinct reconstruction networks, leading to disparate distributions across its dimensions. We apply a BatchNorm [17] layer (without Affine) to $E(A_n)$ and then perform reconstruction residuals selection to ensure a consistent distribution across the dimensions of $E(A_n)$.

The discriminator is implemented using a basic MLP with upsampling layers to map anomaly scores from feature resolution to image resolution. During the training phase of RealNet, we do not use any data augmentation for the synthesis of anomalous images, and maintain an equal ratio between normal images and synthetic anomalous images. In the process of image blending, we uniformly sample the opacity δ from 0.5 to 1.0 in Eq. (3). The training of RealNet is performed on a single NVIDIA GeForce RTX 3090, with an approximate average training time of 2 hours.

C. More results

C.1. Experimental results on BTAD

We evaluate the anomaly detection and localization performance of RealNet and alternative methods on the BTAD dataset [24], with the results shown in Tab. S1. Although SIA does not show a significant performance improvement compared to DTD [5] due to the absence of complex structural anomalies in the three industrial products of the BTAD dataset [24], RealNet demonstrates state-of-the-art performance in anomaly detection and localization when compared to other methods, without any structural or hyperparameter tuning.

Algorithm S1 SDAS with DDIM [34]

Input: diffusion model $\epsilon_\theta(x_t, t)$, perturbation variance Σ , anomaly strength s
 $x_T \sim \mathcal{N}(0, \mathbf{I})$
for all t from T to 1 **do**
 $x_{t-1} \sim \mathcal{N}(\sqrt{\bar{\alpha}_{t-1}}(\frac{x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(x_t, t)}{\sqrt{\bar{\alpha}_t}}) + \sqrt{1 - \bar{\alpha}_{t-1}} \epsilon_\theta(x_t, t), s \Sigma)$
end for
return x_0

Table S1. Comparison of RealNet with alternative anomaly detection methods on the BTAD dataset [24], employing Image AUROC (%) and Pixel AUROC (%) as evaluation metrics.

Category	VT-ADL [24]	P-SVDD [43]	FastFlow [46]	SPADE [6]	RD++ [38]	RealNet (SIA)	RealNet (DTD [5])
01	(-, 99)	(95.7, 91.6)	(-, 95)	(91.4, 97.3)	(96.8, 96.2)	(100.0 , 98.2)	(100.0 , 98.1)
02	(-, 94)	(72.1, 93.6)	(-, 96)	(71.4, 94.4)	(90.1 , 96.4)	(88.6, 96.3)	(87.5, 96.3)
03	(-, 77)	(82.1, 91.0)	(-, 99)	(99.9, 99.1)	(100.0 , 99.7)	(99.6, 99.4)	(99.4, 99.6)
AVG	(-, 90.0)	(83.3, 92.1)	(-, 96.7)	(87.6, 96.9)	(95.6, 97.4)	(96.1 , 97.9)	(95.7, 98.0)

Table S2. Comparison of RealNet with alternative anomaly detection methods on the VisA dataset [55], employing Image AUROC (%) and Pixel AUROC (%) as evaluation metrics.

Category	SPADE [6]	FastFlow [46]	DRAEM [48]	PatchCore [30]	RealNet (SIA)	RealNet (DTD [5])
Candle	(91.0, 97.9)	(92.8, 94.9)	(91.8, 96.6)	(98.6 , 99.5)	(96.1, 99.1)	(95.0, 99.0)
Capsules	(61.4, 60.7)	(71.2, 75.3)	(74.7, 98.5)	(81.6, 99.5)	(93.2 , 98.7)	(88.1, 97.6)
Cashew	(97.8 , 86.4)	(91.0, 91.4)	(95.1, 83.5)	(97.3, 98.9)	(97.8 , 98.3)	(95.9, 97.6)
Chewing gum	(85.8, 98.6)	(91.4, 98.6)	(94.8, 96.8)	(99.1, 99.1)	(99.9, 99.8)	(100.0 , 99.8)
Fryum	(88.6, 96.7)	(88.6, 97.3)	(97.4 , 87.2)	(96.2, 93.8)	(97.1, 96.2)	(95.3, 95.2)
Macaroni1	(95.2, 96.2)	(98.3, 97.3)	(97.2, 99.9)	(97.5, 99.8)	(99.8 , 99.9)	(98.2, 99.7)
Macaroni2	(87.9, 87.5)	(86.3, 89.2)	(85.0, 99.2)	(78.1, 99.1)	(95.2 , 99.6)	(91.8, 99.3)
PCB1	(72.1, 66.9)	(77.4, 75.2)	(47.6, 88.7)	(98.5 , 99.9)	(98.5 , 99.7)	(97.1, 99.4)
PCB2	(50.7, 71.1)	(61.9, 67.3)	(89.8, 91.3)	(97.3, 99.0)	(97.6 , 98.0)	(97.5, 97.8)
PCB3	(90.5, 95.1)	(74.3, 94.8)	(92.0, 98.0)	(97.9, 99.2)	(99.1 , 98.8)	(97.6, 98.4)
PCB4	(83.1, 89.0)	(80.9, 89.9)	(98.6, 96.8)	(99.6, 98.6)	(99.7 , 98.6)	(99.2, 98.6)
Pipe fryum	(81.1, 81.8)	(72.0, 87.3)	(100.0 , 85.8)	(99.8, 99.1)	(99.9, 99.2)	(99.9, 98.6)
AVG	(82.1, 85.6)	(82.2, 88.2)	(88.7, 93.5)	(95.1, 98.8)	(97.8 , 98.8)	(96.3, 98.4)

C.2. Experimental results on VisA

We present the performance of RealNet and alternative methods on the VisA dataset under the one-class protocol [55] in Tab. S2. RealNet achieves the best performance in both anomaly detection and localization. Compared to DTD [5], the RealNet trained using SIA shows an improvement of 1.5% in Image AUROC and 0.4% in Pixel AUROC.

C.3. Supplementary ablation studies

To further investigate RealNet’s anomaly detection performance on the MVTec-AD dataset [3], we examine various backbones and reconstruction feature dimension settings. As shown in Tab. S3, when WideResNet50 [47] is employed as the backbone and the reconstruction feature dimensions $\{m_1, \dots, m_K\}$ are reduced from $\{256, 512, 512, 256\}$ to $\{128, 256, 256, 128\}$, there is a slight decrease of 0.16% in Image AUROC. Despite this reduction, RealNet maintains its competitive performance compared to other

methods. Additionally, the adoption of EfficientNetB4 [35] and ResNet34 [13] as backbones also results in competitive performance, demonstrating the effectiveness of RealNet across various settings.

C.4. Computational efficiency analysis

We investigate the computational efficiency and detection performance of three different multi-scale feature reconstruction architectures on the MVTec-AD dataset [3], as illustrated in Fig. S2. To provide a comprehensive analysis, Tab. S4 presents the inference speed, model size (including backbone), and anomaly detection performance of these architectures. The inference is performed using a single Nvidia GeForce RTX 3090, with all other settings adhering to the specifications detailed in Sec. 4.1.

We utilize a consistent reconstruction network based on the U-Net model with skip connections across three distinct architectures. The employed U-Net model initiates with a

Table S3. Performance evaluation of RealNet with varying backbones and reconstruction feature dimension settings on the MVTec-AD dataset [3], employing Image AUROC (%), Pixel AUROC (%), and PRO (%) as evaluation metrics.

Backbone	EfficientNetB4 [35]	ResNet34 [13]	WideResNet50 [47]	
$\{m_1, \dots, m_K\}$	$\{24, 32, 56, 160\}$	$\{64, 128, 256, 128\}$	$\{128, 256, 256, 128\}$	$\{256, 512, 512, 256\}$
Bottle	(100.0 , 98.83, 95.96)	(100.0 , 98.56, 95.91)	(100.0 , 99.41 , 94.37)	(100.0 , 99.30, 95.62)
Cable	(96.36, 96.33, 88.61)	(96.31, 96.32, 88.68)	(98.35, 98.01, 92.99)	(99.19 , 98.10 , 93.38)
Capsule	(97.97, 99.16, 91.46)	(96.81, 98.78, 87.87)	(99.44, 99.39 , 79.76)	(99.56 , 99.32, 84.48)
Carpet	(100.0 , 98.27, 96.35)	(99.76, 98.37, 94.45)	(99.80, 98.91, 96.32)	(99.84, 99.19 , 96.41)
Grid	(99.92, 99.31, 97.35)	(100.0 , 99.26, 97.39)	(100.0 , 99.55 , 96.38)	(100.0 , 99.51, 97.28)
Hazelnut	(99.89, 98.45, 94.98)	(99.93, 99.35, 94.36)	(100.0 , 99.67, 93.06)	(100.0 , 99.68 , 93.14)
Leather	(100.0 , 99.34, 97.75)	(99.97, 99.40, 98.28)	(100.0 , 99.81 , 96.99)	(100.0 , 99.76, 96.22)
Metal Nut	(99.07, 96.90, 92.65)	(99.17, 96.68, 93.34)	(99.90 , 98.75 , 95.10)	(99.76, 98.58, 94.39)
Pill	(96.10, 94.86, 86.60)	(97.55, 98.23, 93.17)	(97.85, 99.19 , 80.73)	(99.13 , 99.02, 91.04)
Screw	(92.95, 99.05, 92.68)	(96.99, 99.09, 89.57)	(97.99, 99.28, 88.60)	(98.83 , 99.45 , 87.90)
Tile	(99.49, 95.69, 92.10)	(99.93, 97.40, 91.65)	(100.0 , 99.27, 97.20)	(99.96, 99.44 , 97.70)
Toothbrush	(99.44, 98.90, 92.39)	(100.0 , 98.26, 91.74)	(100.0 , 99.26 , 91.22)	(99.44, 98.71, 91.57)
Transistor	(99.58, 98.57 , 93.63)	(99.33, 97.70, 88.53)	(99.79, 98.26, 83.34)	(100.0 , 98.00, 92.92)
Wood	(98.77, 94.47, 92.67)	(98.16, 96.35, 91.46)	(99.56 , 98.22 , 90.76)	(99.21, 98.22 , 90.54)
Zipper	(99.71, 98.01, 91.68)	(99.90 , 98.55, 93.91)	(99.74, 99.20 , 90.73)	(99.82, 99.17, 93.43)
AVG	(98.62, 97.74, 93.12)	(98.92, 98.15, 92.69)	(99.49, 99.07 , 91.17)	(99.65 , 99.03, 93.07)

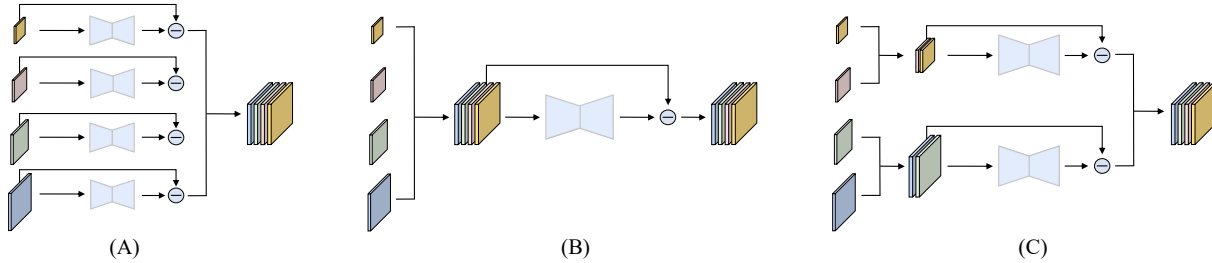


Figure S2. Various architectures of multi-scale feature reconstruction for anomaly detection. (A) Independent Reconstruction Architecture uses separate networks for multi-scale feature reconstruction. (B) Fully Aligned Feature Reconstruction Architecture aligns all features for reconstruction. (C) Neighboring Aligned Feature Reconstruction Architecture aligns and reconstructs neighboring resolution features.

stack of residual layers and down-sampling layers, gradually decreasing the spatial dimensions while increasing the number of channels. Subsequently, the model utilizes a stack of residual layers and up-sampling layers to inversely reconstruct features. Throughout this process, skip connections are incorporated at equivalent spatial resolutions to ensure a smooth and logical flow.

Specifically, architecture **A** adopts separate reconstruction networks to reconstruct multi-scale features without the need for feature interpolation or alignment. This method ensures outstanding anomaly detection performance while maintaining high computational efficiency. With a resolution of 256×256 and reconstruction feature dimensions of $\{256, 512, 512, 256\}$, architecture **A** with model size of 2.2 GB achieves a rapid inference speed of 31.93 FPS. And it can perform inference using only 4GB of GPU memory.

Concurrently, it attains an Image AUROC of 99.65% and a Pixel AUROC of 99.03%. By decreasing the reconstruction feature dimensions to $\{128, 256, 256, 128\}$, architecture **A** reduces the model size to 0.74 GB and achieves a higher inference speed of 40.42 FPS, while preserving an Image AUROC of 99.49% and a Pixel AUROC of 99.07%. Furthermore, at a high resolution of 512×512 , it delivers an inference speed of 13.53 FPS, along with an Image AUROC of 99.40% and a Pixel AUROC of 98.71%. These inference speeds indicate that architecture **A** satisfies the real-time requirements for industrial inspection applications.

Regarding architecture **B**, as referenced in [33, 37, 44], it is used to align the multi-scale features of a small pre-trained network. As aligning down-sampled features will reduce the resolution of model detection and cause predictable performance loss, the experiment only discusses

Table S4. Performance evaluation of various reconstruction architectures on the MVTec-AD dataset [3]. The metrics include Image AUROC (%), Pixel AUROC (%), and PRO (%).

	Speed (FPS) \uparrow	Model Size (GB) \downarrow	Metrics \uparrow
$\{m_1, \dots, m_K\}$ is $\{128, 256, 256, 128\}$ and image size is 256×256			
A	40.42	0.74	(99.49, 99.07 , 91.17)
$\{m_1, \dots, m_K\}$ is $\{256, 512, 512, 256\}$ and image size is 256×256			
A	31.93	2.20	(99.65 , 99.03, 93.07)
B	10.83	7.22	(98.44, 98.17, 94.27)
C	22.39	3.75	(99.62, 98.90, 94.71)
$\{m_1, \dots, m_K\}$ is $\{256, 512, 512, 256\}$ and image size is 512×512			
A	13.53	2.20	(99.40, 98.71, 94.01)

up-sampling alignment. Compared to architecture **A**, architecture **B** reconstructs the interpolated features, significantly reducing computational efficiency and increasing model size. Moreover, due to the limited number of normal images, the overly large reconstruction network in architecture **B** is prone to overfitting, resulting in reduced detection performance. Consequently, for large-scale pre-trained networks with high-dimensional features, aligning and reconstructing all features is suboptimal.

Moreover, we observe that utilizing multiple reconstruction networks for feature reconstruction in architecture **A** causes minor deviations in localizing small-area anomalies, resulting in a reduced PRO. To address this, we propose architecture **C**, which aligns and reconstructs features from two neighboring resolution, thereby reducing the number of reconstruction networks, controlling the model size, and striking a balance between computational efficiency and localization accuracy. At a 256×256 resolution, with reconstruction feature dimensions of $\{256, 512, 512, 256\}$, architecture **C** has a 3.75 GB model size and achieves an inference speed of 22.39 FPS, while attaining an Image AUROC of 99.62%, a Pixel AUROC of 98.90%, and a PRO of 94.71%.

In summary, the design of RealNet balances both anomaly detection performance and computational efficiency. The introduction of AFS allows us to flexibly customize models of various sizes to accommodate different usage scenarios. Furthermore, among our three key innovations, both AFS and RRS introduce no additional learnable parameters, ensuring strong interpretability. As for SDAS, it only introduces perturbation during the reverse diffusion process, without requiring any prior knowledge about the distribution of real anomaly images.

C.5. Anomaly detection in multi-class setting

In the multi-class setting [44, 54], anomaly detection is performed across multiple target classes concurrently, without access to sample class labels during both training and inference phases. Learning the data distributions of multiple classes jointly makes the reconstruction more complex.

Table S5. Comparison of RealNet with alternative methods in multi-class anomaly detection on the MVTec-AD dataset [3].

Methods	Image AUROC	Pixel AUROC
DRAEM [48]	88.1	87.2
PaDiM [7]	84.2	89.5
UniAD [44]	96.5	96.8
OmniAL [54]	97.2	98.3
RealNet	97.3	98.4

Table S6. Image quality comparison of SIA with alternative anomaly synthesis approaches on the MVTec-AD dataset [3].

Methods	FID [15] \downarrow	LPIPS [51] \uparrow
DTD [5]	120.52 ± 0.63	0.16 ± 0.00
CutPaste [20]	77.34 ± 0.09	0.11 ± 0.00
NSA [32]	68.76 ± 0.16	0.09 ± 0.01
SIA	60.39 ± 1.26	0.18 ± 0.01

In such settings, previous reconstruction methods tend to output copies of the input images instead of performing selective reconstruction, which leads to a significant decrease in performance. We evaluate the performance of RealNet in multi-class anomaly detection on the MVTec-AD dataset [3] and compare it with alternative state-of-the-art methods. We use DTD [5] for anomaly synthesis as class labels are unavailable during training. The remaining settings are consistent with Sec. 4.1.

The results are shown in Tab. S5. When detecting anomalies across 15 categories of the MVTec-AD dataset [3] concurrently, RealNet achieves an Image AUROC of 97.3% and a Pixel AUROC of 98.4% using a ResNet50 [13] pre-trained on ImageNet [9], surpassing state-of-the-art multi-class anomaly detection methods [44, 54]. To ensure that normal regions can be reconstructed correctly, we do not explicitly constrain the generalization ability of the reconstructed network in RealNet. Instead, we implicitly constrain the reconstruction network to ensure that anomalous regions can be correctly detected by discarding a part of the reconstruction residuals.

C.6. Synthetic anomaly image quality assessment

In this section, we evaluate the quality of anomaly images generated by various anomaly synthesis methods on the MVTec-AD dataset [3]. Specifically, we use the following evaluation metrics:

- FID (Fréchet Inception Distance) [15]: FID measures the distance between the distribution of synthetic anomaly images and real anomaly images, evaluating both the realism and diversity of the synthetic anomaly images. A lower value indicates better performance.

- LPIPS (Learned Perceptual Image Patch Similarity) [51]: We employ cluster-based LPIPS [11] to evaluate the diversity of synthetic anomaly images. Supposing a category contains N real anomaly images, we partition the synthesized anomaly images into N groups by finding the lowest LPIPS, then we compute the mean pairwise LPIPS within each group and compute the average of all groups. A higher cluster LPIPS indicates greater diversity.

We employ various anomaly synthesis methods to generate 1,000 anomaly images for evaluation, with each method independently assessed three times. The experimental results are shown in Tab. S6. In comparison to other anomaly synthesis methods, SIA achieves the best FID and LPIPS metrics, highlighting the outstanding performance of SDAS in generating both realistic and diverse anomaly images, and demonstrating the effectiveness of SDAS in improving anomaly detection performance.

D. Visualization

We conduct a comprehensive visual analysis of RealNet on the four datasets. Fig. S3 shows the qualitative results of RealNet in anomaly localization, showcasing its outstanding performance in pixel-level anomaly localization. Figs. S4 and S5 display the anomaly images and normal images generated by SDAS, respectively. Fig. S6 illustrates images synthesized using SIA with localized anomalous regions. Fig. S7 provides an intuitive explanation of pre-training bias, indicating that not all feature maps contribute equally to anomaly detection and localization, which validates the efficacy of AFS.

E. Limitations

In some categories with more texture anomalies, such as the texture categories in MVTec-AD dataset [3], SIA’s performance may slightly underperform when compared to DTD [5]. Given that DTD dataset [5] includes a diverse range of real-world texture images, it effectively simulates common anomaly types in the textural category, such as color, oil, and glue. Nonetheless, SIA excels in the majority of scenarios, outperforming DTD [5] and offering superior capability in synthesizing anomalies in images with intricate structures.

Compared to anomaly synthesis methods based on data augmentation [20, 32] or external data [48], SDAS increases additional offline training time. For instance, we generate 10,000 anomaly images at a resolution of 256×256 for each category, and it will take 6 hours using a single NVIDIA GeForce RTX 3090. However, it is pivotal to clarify that RealNet omits SDAS without any additional computational cost during inference and real-world applications. Therefore, we believe that the slight increase in training time to enhance performance is necessary and worthwhile.

In order to achieve higher computational efficiency, we do not upsample multi-scale features. Instead, we employ multiple reconstruction networks for feature reconstruction, which reduce the resolution of anomaly detection. The lower feature reconstruction resolution may introduce minor deviations in localizing small anomalous areas, leading to a decrease in PRO. However, we found that increasing the resolution of anomaly detection by reducing the number of reconstruction networks can improve PRO. For instance, architecture C in Fig. S2 achieved a higher PRO score of 94.71%. Furthermore, increasing the resolution of images can also lead to an improvement in PRO, as detailed in Tab. S4.

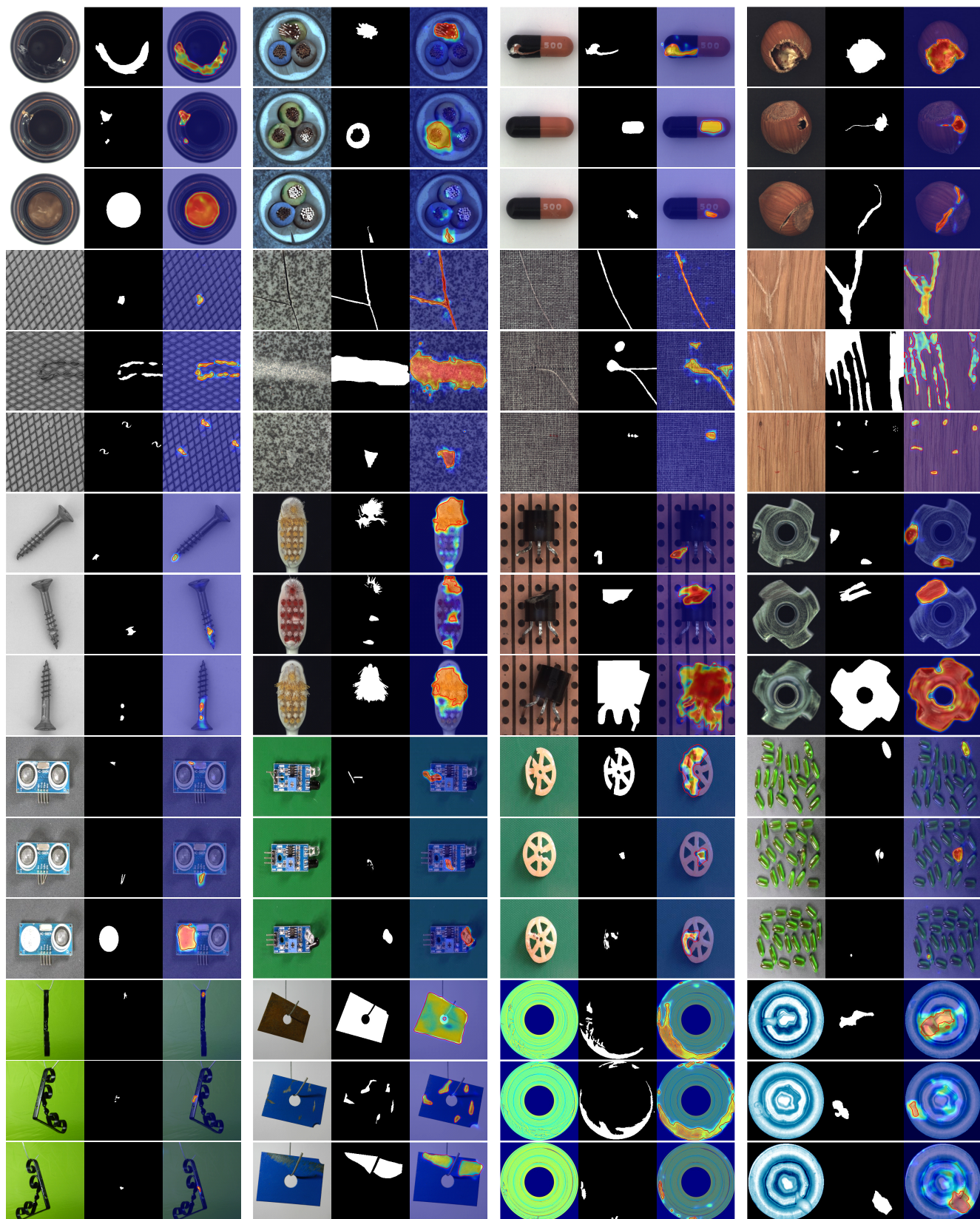


Figure S3. Qualitative results of RealNet. Within each group, from left to right, are the anomaly image, ground-truth, and predicted anomaly score. The examples are from the MVTec-AD [3], MPDD [18], BTAD [24], and VisA [55] datasets.



Figure S4. Anomaly images generated by SDAS. The examples are from the MVtec-AD [3], MPDD [18], BTAD [24], and VisA [55] datasets. Within each group, from top to bottom, the anomaly strength gradually increases.



Figure S5. Normal images generated by SDAS (when $s = 0$). The examples are from the MVTec-AD [3], MPDD [18], BTAD [24], and VisA [55] datasets.

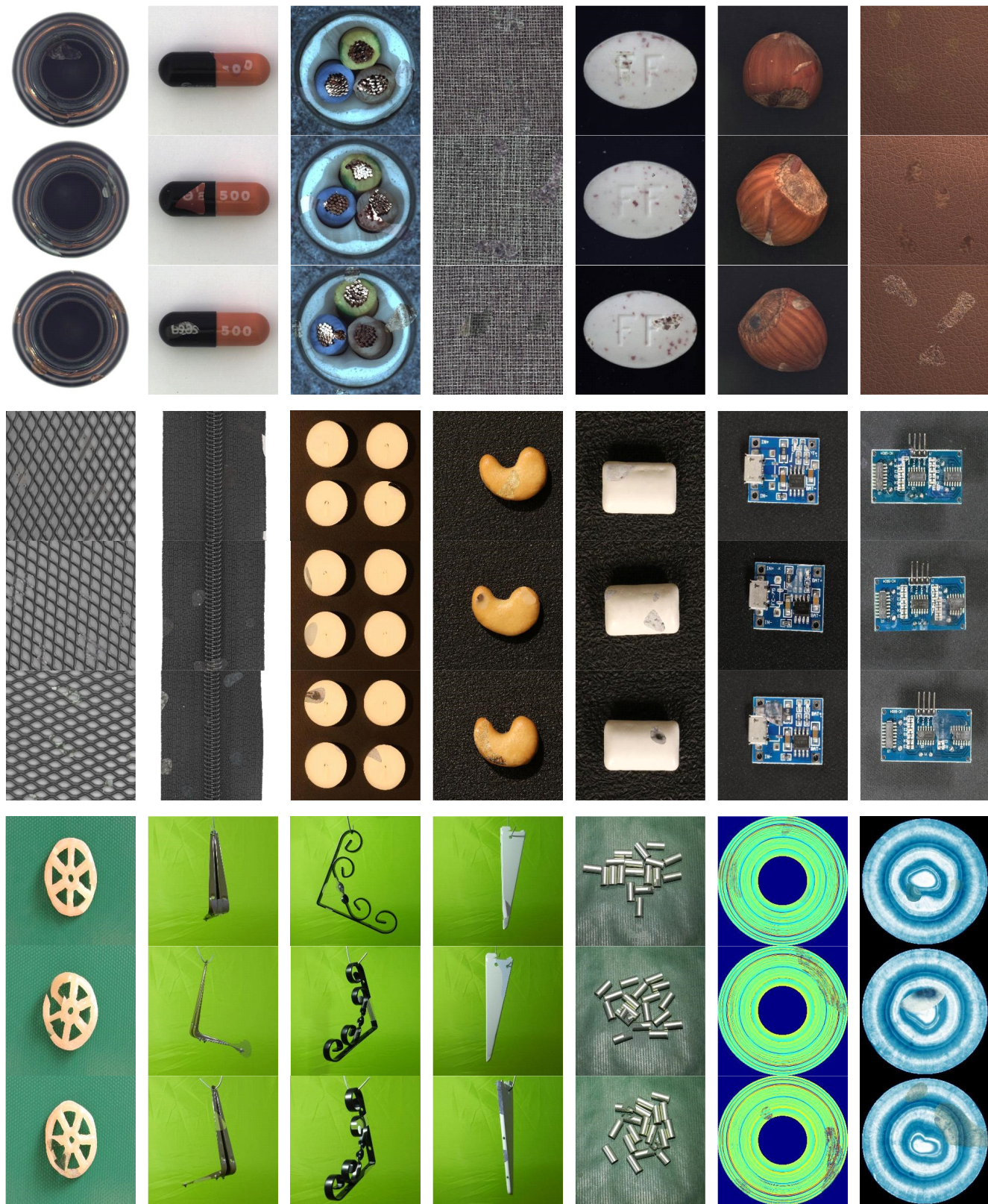


Figure S6. Local anomaly images synthesized by SIA. The examples are from the MVTec-AD [3], MPDD [18], BTAD [24], and VisA [55] datasets. Within each group, from top to bottom, the anomaly strength gradually increases.

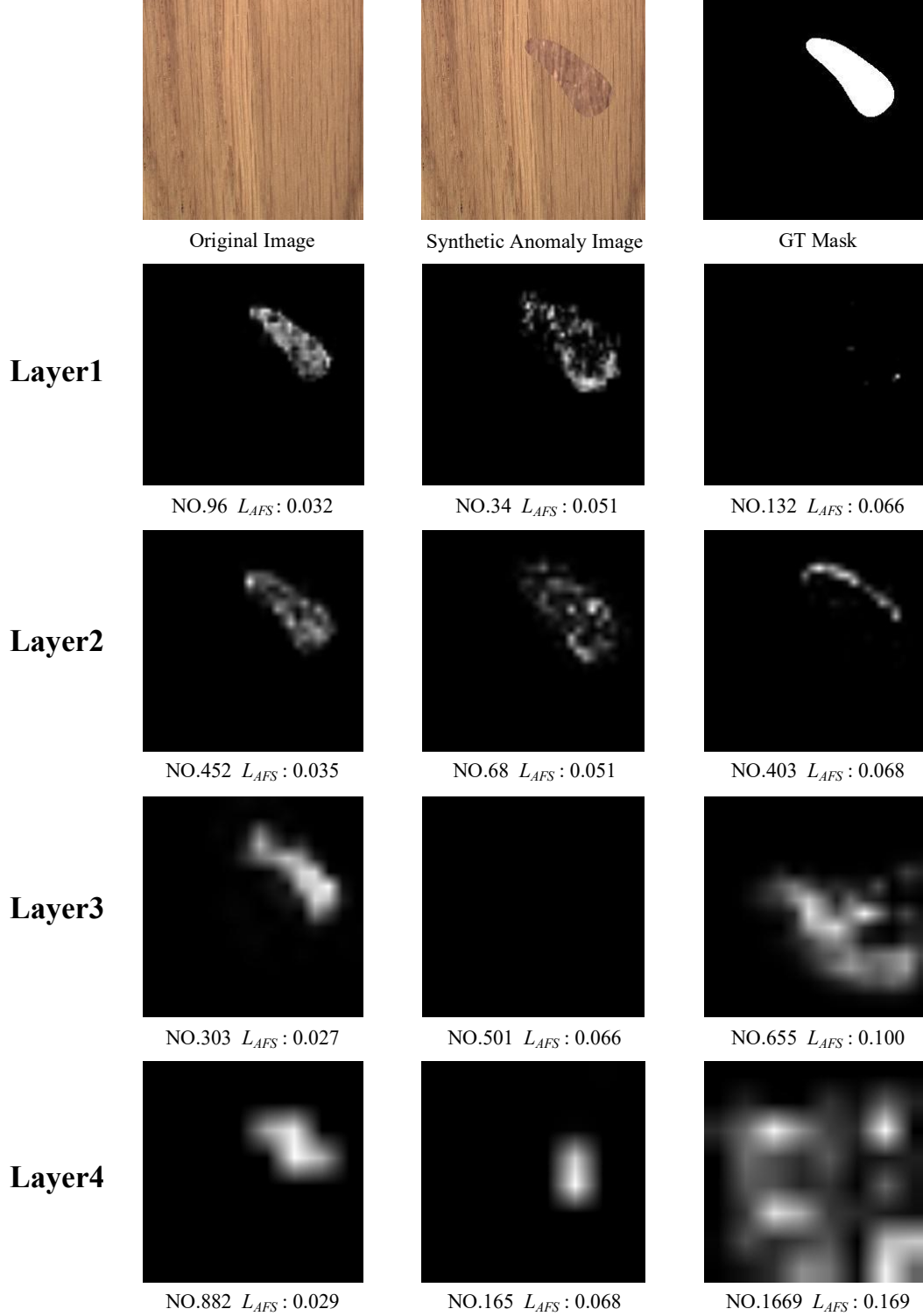


Figure S7. Visualization of AFS. For an original image and a synthetic anomaly image, we visualize the normalized difference between their corresponding feature maps across different layers of a pre-trained WideResNet50 [47]. From top to bottom, the feature map respectively come from the first layer to the fourth layer. Each feature map is labelled with its index in the layer and the corresponding AFS loss. From left to right, the localization performance of the feature maps gradually decreases. Our visualization intuitively demonstrates the localization bias caused by pre-training, indicating that not all feature maps contribute equally to anomaly detection and localization, as well as emphasizing the effectiveness of AFS.