

# SeaS: Few-shot Industrial Anomaly Image Generation with Separation and Sharing Fine-tuning

Zhewei Dai<sup>1,\*</sup>, Shilei Zeng<sup>1,\*</sup>, Haotian Liu<sup>1</sup>, Xurui Li<sup>1</sup>, Feng Xue<sup>4</sup>, Yu Zhou<sup>1,2,3,†</sup>

<sup>1</sup> School of Electronic Information and Communications, Huazhong University of Science and Technology

<sup>2</sup> Hubei Key Laboratory of Smart Internet Technology, Huazhong University of Science and Technology

<sup>3</sup> Artificial Intelligence Research Institute, Wuhan JingCe Electronic Group Co.,LTD

<sup>4</sup> Department of Information Engineering and Computer Science, University of Trento

{zwdai, shlzeng, htliu.master, xrli.plus, yuzhou}@hust.edu.cn, xuefengbupt@gmail.com

## Abstract

We introduce SeaS, a unified industrial generative model for automatically creating diverse anomalies, authentic normal products, and precise anomaly masks. While extensive research exists, most efforts either focus on specific tasks, i.e., anomalies or normal products only, or require separate models for each anomaly type. Consequently, prior methods either offer limited generative capability or depend on a vast array of anomaly-specific models. We demonstrate that U-Net’s differentiated learning ability captures the distinct visual traits of slightly-varied normal products and diverse anomalies, enabling us to construct a unified model for all tasks. Specifically, we first introduce an Unbalanced Abnormal (UA) Text Prompt, comprising one normal token and multiple anomaly tokens. More importantly, our Decoupled Anomaly Alignment (DA) loss decouples anomaly attributes and binds them to distinct anomaly tokens of UA, enabling SeaS to create unseen anomalies by recombining these attributes. Furthermore, our Normal-image Alignment (NA) loss aligns the normal token to normal patterns, making generated normal products globally consistent and locally varied. Finally, SeaS produces accurate anomaly masks by fusing discriminative U-Net features with high-resolution VAE features. SeaS sets a new benchmark for industrial generation, significantly enhancing downstream applications, with average improvements of +8.66% pixel-level AP for synthesis-based AD approaches, +1.10% image-level AP for unsupervised AD methods, and +12.79% IoU for supervised segmentation models. Code is available at <https://github.com/HUST-SLOW/SeaS>.

## 1. Introduction

In the industrial scenario, generative models are used to synthesise various visual elements, which meet the require-

\* Contributed Equally, † Corresponding Authors.

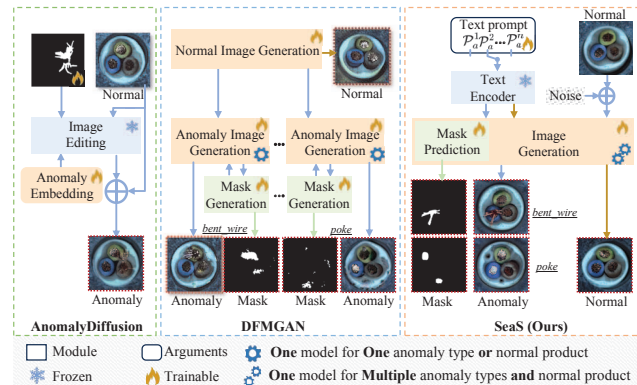


Figure 1. (a) **AnomalyDiffusion** only generates anomalies and edits them onto the input normal images guided by the predefined masks. (b) **DFMGAN** trains multiple dedicated generators per anomaly type or normal product, and it cannot produce accurate masks. (c) **SeaS** trains a unified model capable of generating anomaly images and masks for multiple anomaly types, as well as normal images.

ments of different anomaly detection (AD) methods and supervised segmentation models as below.

- Generating pseudo-anomalies for synthesis-based AD approaches [9, 41].
- Generating pseudo-normal images for unsupervised AD methods [15, 26, 32].
- Generating complete anomaly images and corresponding accurate masks for training supervised segmentation models.

These requirements above have been covered by previous algorithms. CutPaste [21] creates anomalies by pasting cropped normal regions onto normal product images (short for normal images). AnomalyDiffusion [17] generates anomalies using diffusion models and edits anomalies onto normal images guided by pre-defined masks (Fig. 1(a)), but it cannot create pseudo-normal images and might suffer from misaligned masks. DFMGAN [11] trains separate models for each normal product and anomaly type

(Fig. 1(b)) but cannot produce accurate masks, limiting its effectiveness in training supervised segmentation models. In summary, existing methods either focus only on the generation of normal products or anomalies, or require multiple isolated models to complete all tasks separately. They cannot flexibly use a unified model to tackle them all, i.e., achieving diverse anomalies, authentic normal products, and pixel-accurate masks. In this paper, we study the unified industrial anomaly generation solution, meeting the needs of various downstream tasks.

The novelty of this work stems from a key observation on a single industrial production line: *normal products exhibit a globally consistent surface with minor local variations, whereas anomalies exhibit high variability*. These characteristics can be effectively captured by U-Net due to its differential learning capability in a diffusion process. Building on this insight, we propose a **Separation and Sharing Fine-tuning method (SeaS)**, using a shared U-Net to model distinct variations. Firstly, to explicitly model the variations of normal products and anomalies, we propose **Unbalanced Abnormal (UA) Text Prompt**. Its unbalanced design includes one normal token and multiple anomaly tokens, thus decoupling the presentations of the slightly varied normal product surface and diverse anomaly semantics. Secondly, to learn highly-diverse anomalies, we propose a **Decoupled Anomaly Alignment (DA) loss** to bind the attributes of anomalies to different anomaly tokens of UA. Recombining the decoupled attributes may produce anomalies that have never been seen in the training dataset, therefore increasing the diversity of generated anomalies. Thirdly, for slightly varied normal products, we propose the **Normal-image Alignment (NA) loss**. It enables the network to learn the key features of the normal product from normal images, so that the normal token of UA expresses the products’ global consistency. The two training processes above are separated but conducted on a shared U-Net. SeaS enables U-Net to simultaneously model the different variations in both normal products and anomalies, representing the discriminative features for mask prediction. However, the low-resolution features of U-Net lead to a coarse mask when predicted directly. Thus, we propose a **Refined Mask Prediction (RMP) branch**. It combines U-Net features with high-resolution VAE features to generate accurate and crisp masks progressively. The generated anomaly images achieve IS scores by 1.88 (MVTec AD), 1.27 (VisA), and 1.95 (MVTec 3D AD), with IC-LPIPS of 0.34, 0.26, and 0.30. On multiple datasets, SeaS-generated images boost synthesis-based AD approaches by an average +8.66% pixel-level AP, improve unsupervised AD methods by an average +1.10% image-level AP, and enhance supervised segmentation models by an average +12.79% IoU.

In summary, the key contribution of our approach lies in:

- We propose a unified generative model for industrial

visual elements. It achieves diverse anomalies, globally consistent normal products, and pixel-level accurate masks using only one model, which sets a new standard for this field.

- The newly designed separated and shared fine-tuning models different variations of normal products and anomalies, enabling precise control over their generation, and obtaining discriminative features for mask prediction.
- SeaS greatly improves the performance of various synthesis-based and unsupervised AD methods, and empowers supervised segmentation models with decent performance.

## 2. Related Work

**Anomaly Image Generation.** Early non-generative methods [10, 21, 41] rely on data augmentation to synthesize pseudo-anomalies, but suffer from low fidelity due to inconsistent anomaly patterns. Some generative methods [13, 17, 34] only generate anomalies and merge them into the real normal images. NSA [34] uses Poisson Image Editing [30] to fuse the cropped normal region. However, these methods cannot create pseudo-normal images and require anomaly masks as inputs, with unreasonable mask positions compromising fidelity and consistency. GAN-based methods [11, 28, 42] generate the entire anomaly images. DFMGAN[11] trains multiple isolated models to generate normal images and anomaly images for each anomaly type, and the produced masks often do not align accurately with anomalies, limiting their utility in training supervised segmentation models. Different from these existing approaches, we propose a unified generative model based on Stable Diffusion to generate diverse anomalies, globally consistent normal products and pixel-level accurate masks.

**Fine-tuning Diffusion Models.** Fine-tuning is a potent strategy for enhancing specific capabilities of pre-trained diffusion models [6, 12, 43]. Personalized methods [8, 12, 33] utilize a small set of images to fine-tune the diffusion model, thereby generating images of the same object. Several methods for multi-concept image fine-tuning [1, 14, 19, 20, 37] use cross-attention maps to align embeddings with individual concepts in the image. Nevertheless, they do not consider the different variations in different image regions, which is important for industrial anomaly image generation. Thus, we propose a separation and sharing fine-tuning strategy to model the different degrees of variations of anomalies and normal products, which independently learns products and anomalies on a shared U-Net.

**Mask Prediction with Generation Method.** Previous methods on mask prediction for generated images are mainly based on features in GANs [22, 44]. However, these approaches do not guarantee the generation of accurate masks for exceedingly small datasets. Based on Stable Diffusion [31], some recent methods, i.e., Diffu-

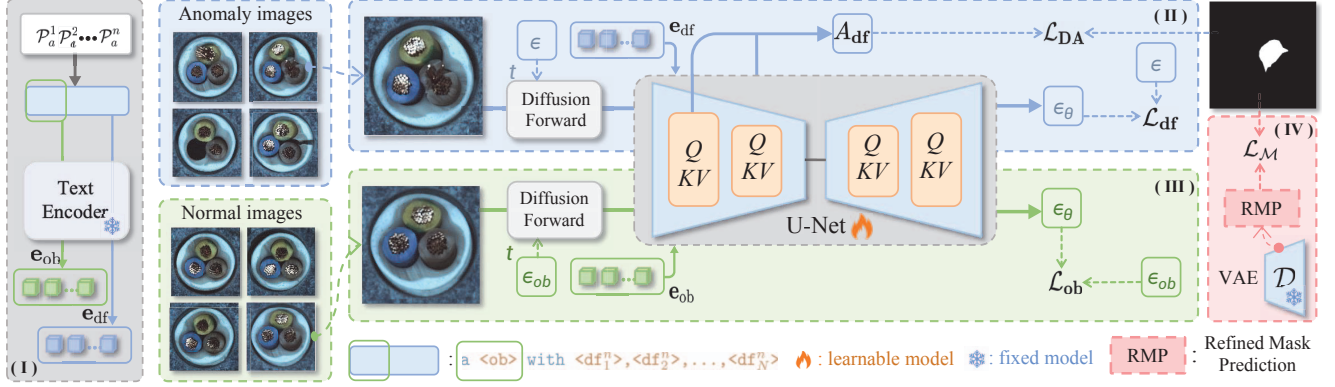


Figure 2. **Overall framework of SeaS.** It consists of four parts: (I) the Unbalanced Abnormal Text Prompt, (II) the Decoupled Anomaly Alignment for aligning the anomaly tokens  $\langle df_n \rangle$  to the anomaly area of abnormal images, (III) the Normal-image Alignment for maintaining authenticity through normal images, and (IV) the Refined Mask Prediction branch for generating accurate masks.

Mask [36], DatasetDM [35] and DatasetDiffusion [27], produce masks by exploiting the potential of the cross-attention maps. However, due to the low resolution of the cross-attention maps, they are directly interpolated to a higher resolution to match the image size without any auxiliary information, which leads to significant boundary uncertainty. We incorporate the high-resolution features from the VAE decoder as auxiliary information for resolution retrieving, fusing them with U-Net decoder features, which are discriminative due to the modelling of the different variations in normal products and anomalies, to generate accurate high-resolution masks.

### 3. Method

The training phase of the proposed Separation and Sharing (SeaS) Fine-tuning strategy is shown in Fig. 2. In Sec. 3.1, we introduce the preliminaries of our approach. In Sec. 3.2, we first design an Unbalanced Abnormal Text Prompt, which contains a set of tokens that characterize normal products and anomalies separately. Subsequently, we propose the Decoupled Anomaly Alignment (DA) loss to bind anomaly image regions to anomaly tokens, and leverage Normal-image Alignment (NA) loss to empower normal token to express the globally-consistent normal product surface. The two training processes are implemented separately for abnormal and normal images but on a shared U-Net architecture. Then, based on the well-trained U-Net, we design a Refined Mask Prediction branch to generate accurate masks corresponding to the generated anomaly images in Sec. 3.3. Finally, we detail the generation of abnormal image-mask pairs and normal images in Sec. 3.4.

#### 3.1. Preliminaries

**Stable Diffusion.** Given an input image  $x_0$ , Stable Diffusion [31] firstly transforms  $x_0$  into a latent space as  $z = \varepsilon(x_0)$ , and then adds a randomly sampled noise  $\epsilon \sim N(0, \mathbf{I})$  into  $z$  as  $\hat{z}_t = \alpha_t z + \beta_t \epsilon$ , where  $t$  is the randomly sampled timestep. Then, the U-Net is employed to predict the noise

$\epsilon$ . Let  $c_\theta(\mathcal{P})$  be the CLIP text encoder that maps conditioning text prompt  $\mathcal{P}$  into a conditioning vector  $\mathbf{e}$ . The training loss of Stable Diffusion can be stated as follows:

$$\mathcal{L}_{SD} = \mathbb{E}_{z=\varepsilon(x_0), \mathcal{P}, \epsilon \sim N(0, \mathbf{I}), t} \left[ \|\epsilon - \epsilon_\theta(\hat{z}_t, t, \mathbf{e})\|_2^2 \right] \quad (1)$$

where  $\epsilon_\theta$  is the predicted noise.

**Cross-Attention Map in U-Net.** Aiming to control the generation process, the conditioning mechanism is implemented by calculating cross-attention between the conditioning vector  $\mathbf{e} \in \mathbb{R}^{Z \times C_1}$  and image features  $\mathbf{v} \in \mathbb{R}^{r \times r \times C_2}$  of the U-Net inner layers [7, 16, 39]. The cross-attention map  $A^{m,l} \in \mathbb{R}^{r \times r \times Z}$  can be calculated as:

$$A^{m,l} = \text{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right), Q = \phi_q(\mathbf{v}), K = \phi_k(\mathbf{e}) \quad (2)$$

where  $Q \in \mathbb{R}^{r \times r \times C}$  denotes a query projected by a linear layer  $\phi_q$  from  $\mathbf{v}$ ,  $r$  is the resolution of the feature map in U-Net, and  $l$  is the index of the U-Net inner layer.  $K \in \mathbb{R}^{Z \times C}$  denotes a key through another linear layer  $\phi_k$  from  $\mathbf{e}$ , and  $Z$  is the number of text embeddings after padding.

#### 3.2. Separation and Sharing Fine-tuning

**Unbalanced Abnormal Text Prompt.** Through the experimental observation, we found that the typical text prompt, like a photo of a bottle with defect [18], or damaged bottle [46], is suboptimal for industrial anomaly generation. The balanced semantic words for normal products and anomalies may fail to capture their differential variation degrees. Therefore, we design the Unbalanced Abnormal (UA) Text Prompt for each anomaly type of each product, i.e.,

$\mathcal{P} = a \langle ob \rangle$  with  $\langle df_1 \rangle, \langle df_2 \rangle, \dots, \langle df_N \rangle$  where  $\langle ob \rangle$  and  $\langle df_n \rangle$  ( $n \in \{1, 2, \dots, N\}$ ) are the tokens of the industrial normal products (short for Normal Token) and the anomalies (short for Anomaly Token) respectively. We use a set of  $N$  Anomaly Tokens for each anomaly type, with different sets corresponding to different

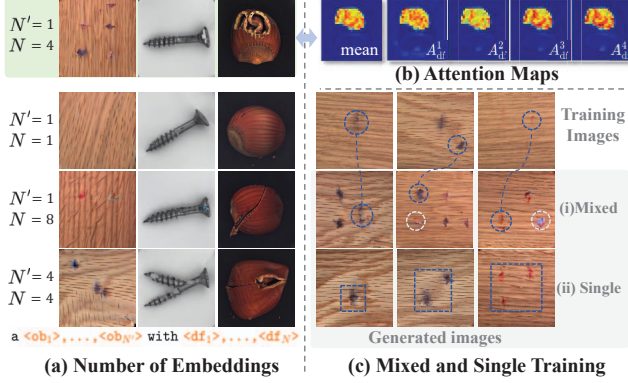


Figure 3. (a) Generated images with the different number of tokens. (b) Cross-attention maps. (c) Examples of diverse generated images.

anomaly types. As shown in Fig. 3, in SeaS, we separately employ normal images to train the embedding corresponding to  $\langle ob \rangle$ , and abnormal images to train the embeddings corresponding to  $\langle df_n \rangle$ . Experimental observations indicate that one  $\langle ob \rangle$  is sufficient to express normal product, while multiple  $\langle df_n \rangle$  are necessary for controlling the generation of anomalies. As shown in Fig. 3(a), when we use the UA prompt  $\mathcal{P}$  (the dotted green box in (a)), the cross-attention maps in (b) show that different tokens have different responses in the abnormal regions, which indicates that they focus on different attributes of the anomalies, and performing the average operation on the cross-attention maps produces never-seen anomalies. When we use only one  $\langle df \rangle$ , it is difficult to align it to several different anomalies that belong to the same category. Therefore, during inference, if the denoised anomaly feature has a larger distance to  $\langle df \rangle$ , it will be assigned a smaller response by the U-Net, which leads to the ‘‘anomaly missing’’ phenomenon, e.g., the generated images in the case of  $(N' = 1, N = 1)$ . In addition, if we utilize a large number of  $\langle df_n \rangle$ , we find that each  $\langle df_n \rangle$  may focus on some local properties of an anomaly, such a case increases the diversity but may reduce the authenticity of the anomalies, as shown in the case  $N' = 1, N = 8$ . Similarly, if we use multiple learnable  $\langle ob \rangle$ , e.g.,  $N' = 4, N = 4$ , each  $\langle ob \rangle$  pays attention to the local character of the normal product, which may reduce the global consistency of the normal product.

**Decoupled Anomaly Alignment.** Given a few abnormal images  $x_{df}$  and their corresponding masks, we aim to align the anomaly tokens  $\langle df_n \rangle$  to the anomaly area of  $x_{df}$  by tuning the U-Net and the learnable embedding corresponding to  $\langle df_n \rangle$ . Therefore, we propose the Decoupled Anomaly Alignment (DA) loss, i.e.,

$$\mathcal{L}_{DA} = \sum_{l=1}^L \left( \left\| \frac{1}{N} \sum_{n=1}^N A_{df}^{n,l} - M^l \right\|^2 + \|A_{ob}^l \odot M^l\|^2 \right) \quad (3)$$

where  $A_{df}^{n,l} \in \mathbb{R}^{r \times r \times 1}$  is the cross-attention map corresponding to the  $n$ -th anomaly token  $\langle df_n \rangle$ ,  $N$  is the num-

ber of anomaly token in  $\mathcal{P}$ .  $L$  is the total number of U-Net layers used in alignment.  $M^l$  is the binary mask with  $r \times r$  resolution, where the abnormal area is 1 and the background is 0.  $A_{ob}^l \in \mathbb{R}^{r \times r \times 1}$  is the cross-attention map corresponding to the normal token  $\langle ob \rangle$ ,  $\odot$  is the element-wise product. DA loss performs the mandatory decoupling of the anomaly and the normal product. The first term of DA loss is to align the abnormal area to  $\langle df_n \rangle$  according to the mask  $M^l$ . The second term of DA loss reduces the response value of  $A_{ob}^l$  in the abnormal area, which prevents  $\langle ob \rangle$  from aligning to the abnormal area of  $x_{df}$ . Further analysis of how the DA loss ensures the diversity of anomalies is provided in Appendix A.2. Therefore, the total loss for the anomaly image  $x_{df}$  is:

$$\mathcal{L}_{df} = \mathcal{L}_{DA} + \|\epsilon_{df} - \epsilon_{\theta}(\hat{z}_{df}, t_{df}, \mathbf{e}_{df})\|_2^2 \quad (4)$$

In second term of Eq. 4, we use random noises  $\epsilon_{df}$  and timesteps  $t_{df}$  to perform forward diffusion on abnormal images  $x_{df}$ , then obtain the noisy latent  $\hat{z}_{df}$ . The conditioning vector  $\mathbf{e}_{df} \in \mathbb{R}^{r \times C_1}$  is used to guide the U-Net in predicting noise, and then calculate the loss with the noise  $\epsilon_{df}$ .

**Normal-image Alignment.** As we discussed, increasing the number of the normal token  $\langle ob \rangle$  leads to a higher diversity, while it may reduce the authenticity of the generated normal image and destruct global consistency. However, aligning only one  $\langle ob \rangle$  to a few of the training images may suffer from the issue of overfitting. Therefore, we add a Normal-image Alignment (NA) loss to overcome such a dilemma, which is stated as follows,

$$\mathcal{L}_{ob} = \|\epsilon_{ob} - \epsilon_{\theta}(\hat{z}_{ob}, t_{ob}, \mathbf{e}_{ob})\|_2^2 \quad (5)$$

Instead of aligning the normal region of  $x_{df}$  to  $\langle ob \rangle$ , in calculating the NA loss, we use random noises  $\epsilon_{ob}$  and timesteps  $t_{ob}$  to perform forward diffusion on the normal images  $x_{ob}$ . Then the noisy latent  $\hat{z}_{ob}$  and the embedding  $\mathbf{e}_{ob}$  corresponding to the normal tokens of  $\mathcal{P}$ , i.e., ‘‘a  $\langle ob \rangle$ ’’, are input into the U-Net in predicting noise, and then calculate the NA loss with  $\epsilon_{ob}$ .

**Mixed Training.** Based on the separated DA loss for abnormal images and NA loss for the normal ones, the objective of Separation and Sharing Fine-tuning is formed as:

$$\mathcal{L} = \mathcal{L}_{df} + \mathcal{L}_{ob} \quad (6)$$

In the training process, instead of training a single U-Net model for each anomaly type, we train a unified U-Net model for each product. Specifically, given a product image set, which contains  $G$  anomaly categories of masked abnormal images and some normal images. We group all the abnormal images of a product into a unified set  $X_{df} = \{x_{df}^1, x_{df}^2, \dots, x_{df}^H\}$ . For each anomaly type, we use  $\mathcal{P}$  with different sets of anomaly tokens. In addition, we sample a fixed number of normal images to consist of the normal training set  $X_{ob} = \{x_{ob}^1, x_{ob}^2, \dots, x_{ob}^P\}$ . During each

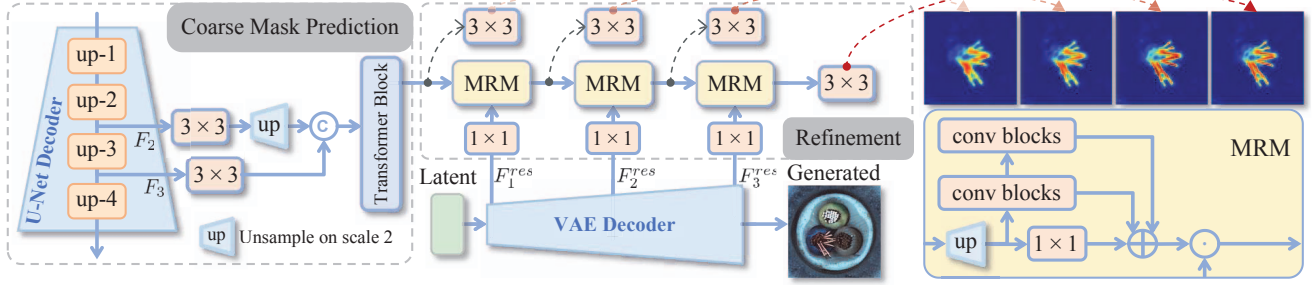


Figure 4. The **Refined Mask Prediction (RMP)** branch during inference. The Coarse Feature Extraction utilizes features from the up-2 and up-3 layers of the U-Net Decoder to extract coarse features. The cascaded Mask Refinement Module (MRM) further obtains the mask accurately aligned with the anomaly with the assistance of high-resolution features of the VAE Decoder.

step of our fine-tuning process, we sample same number of images from both  $X_{df}$  and  $X_{ob}$ , and mixed them into a batch. We found that such a mixed training strategy not only alleviates the overfitting caused by the limited number of each anomaly type, but also increases the diversity of the anomaly image, while still maintaining reasonable authenticity, as shown in Fig. 3(c), (i) indicates that the model with mixed training may generate new anomalies, e.g., the anomalies inside the dotted white line. In contrast, the anomalies in (ii) overfit the training images. More ablation studies on the mixed training strategy are shown in Tab. 12 in appendix A.5.

### 3.3. Refined Mask Prediction

The design of the separated and shared approach enables U-Net to simultaneously model the different degrees of variations in both normal products and anomalies, representing the discriminative features for mask prediction. To further obtain pixel-accurate masks, we design a cascaded Refined Mask Prediction (RMP) branch, which is grafted onto the U-Net trained within SeaS (mentioned in Sec. 3.2). As shown in Fig. 4, RMP consists of two steps, firstly capturing discriminative features from U-Net and secondly combining them with high-resolution features of VAE decoder to generate anomaly-matched masks.

**Coarse Feature Extraction.** The first step aims to extract a coarse but highly-discriminative feature for anomalies from the U-Net decoder. Specifically, let  $F_1 \in \mathbb{R}^{32 \times 32 \times 1280}$  and  $F_2 \in \mathbb{R}^{64 \times 64 \times 640}$  denote the output feature of “up-2” and “up-3” layers of the decoder in U-Net, respectively. We first leverage a  $1 \times 1$  convolution block to compress the channel of  $F_1$  and  $F_2$  to  $\bar{F}_1 \in \mathbb{R}^{32 \times 32 \times 128}$  and  $\bar{F}_2 \in \mathbb{R}^{64 \times 64 \times 64}$ , respectively. Then, we upsample  $\bar{F}_1$  to  $64 \times 64$  resolution and concatenate it with  $\bar{F}_2$ . Finally, four transformer layers are employed to fuse the concatenated features and obtain a unified coarse feature  $\hat{F} \in \mathbb{R}^{64 \times 64 \times 192}$ .

**Mask Refinement Module.** Directly upsampling the coarse feature  $\hat{F}$  to high resolution will result in a loss of anomaly details. Therefore, we design the Mask Refinement Module (MRM) to refine the coarse feature  $\hat{F}$  in a

progressive manner. As shown in Fig. 4, each MRM takes in two features, i.e., the high-resolution features from VAE and the discriminative feature to be refined. Firstly, the discriminative feature is upsampled to align with the high-resolution ones of VAE. To preserve the discriminative ability, the upsampled feature is processed through two chained convolution blocks for capturing multi-scale anomaly features and a  $1 \times 1$  convolution for capturing local features. These features are then summed and multiplied with the VAE features element-wisely to enhance the anomalies’ boundary. Finally, MRM employs a  $3 \times 3$  convolution to fuse the added features and outputs a refined feature.

To refine  $\hat{F}$ , we employ three MRMs positioned in sequence. Each MRM takes the previous MRM’s output as the discriminant feature to be refined, while the first MRM takes  $\hat{F}$  as the discriminative input. For another input of each MRM, we use the outputs from the 1-st, 2-nd, and 3-rd “up-blocks” of the VAE decoder respectively. In this way, the features obtained by the last MRM have the advantages of both high resolution and high discriminability. Finally, we use a  $3 \times 3$  convolution and a softmax to generate the refined anomaly mask  $\hat{M}'_{df} \in \mathbb{R}^{512 \times 512 \times 2}$  using the output of the last MRM.

**Loss Functions.** During training, we use  $x_{df}$  and  $x_{ob}$  as inputs. For  $x_{df}$ , we obtain the coarse mask  $\hat{M}_{df} \in \mathbb{R}^{64 \times 64 \times 2}$  from the Coarse Feature Extraction and  $\hat{M}'_{df}$  after the MRMs. Similarly, for  $x_{ob}$ , we obtain the  $\hat{M}_{ob} \in \mathbb{R}^{64 \times 64 \times 2}$  from Coarse Feature Extraction and directly upsample it to the original resolution, denoted as  $\hat{M}'_{ob} \in \mathbb{R}^{512 \times 512 \times 2}$ . Then we conduct the supervision on both low-resolution and high-resolution predictions as,

$$\mathcal{L}_{\mathcal{M}} = \mathcal{F}(\hat{M}_{df}, \mathbf{M}_{df}) + \mathcal{F}(\hat{M}_{ob}, \mathbf{M}_{ob}) + \mathcal{F}(\hat{M}'_{df}, \mathbf{M}'_{df}) + \mathcal{F}(\hat{M}'_{ob}, \mathbf{M}'_{ob}) \quad (7)$$

where  $\mathcal{F}$  indicates the Focal Loss [24].  $\mathbf{M}_{ob} \in \mathbb{R}^{64 \times 64 \times 1}$  and  $\mathbf{M}'_{ob} \in \mathbb{R}^{512 \times 512 \times 1}$  are used to suppress noise in normal images, with each pixel value set to 0.  $\mathbf{M}_{df} \in \mathbb{R}^{64 \times 64 \times 1}$  and  $\mathbf{M}'_{df} \in \mathbb{R}^{512 \times 512 \times 1}$  are the ground truth masks of abnormal images. More ablation studies on the

effect of normal images in training RMP branch are shown in Tab. 16 and Fig. 11 in appendix A.5.

### 3.4. Inference

During the generation, aiming further to ensure the global consistency of the normal products, we random select a normal image  $x_{ob}$  from  $X_{ob}$  as input, and add random noise to  $x_{ob}$ , which resulting in an initial noisy latent  $\hat{z}_0$ . Next, for the generation of abnormal images,  $\hat{z}_0$  is input into the U-Net for noise prediction, with the process guided by the conditioning vector  $\mathbf{e}_{df}$  (mentioned in Eq. 4), which is corresponding to the whole UA Text Prompt  $\mathcal{P}$ . For generating normal images to further enhance unsupervised AD methods, we use the conditioning vector  $\mathbf{e}_{ob}$  corresponding to the normal tokens of  $\mathcal{P}$ . Regarding the masks corresponding to anomalies, in the final three denoising steps, the RMP branch (Sec. 3.3) leverages the features from the U-Net decoder and VAE decoder to generate the final anomaly mask. Specifically, we average the refined anomaly mask from these steps to obtain the refined mask  $\hat{M}'_{df} \in \mathbb{R}^{512 \times 512 \times 2}$ . Then we take the threshold  $\tau$  for the second channel of  $\hat{M}'_{df}$  to segment the final anomaly mask  $M_{df} \in \mathbb{R}^{512 \times 512 \times 1}$ . The effect of  $\tau$  on the downstream supervised segmentation models is shown in Tab. 18 in appendix A.5. In the last denoising step, the output of the generation model is used as the generated abnormal image.

## 4. Experiments

### 4.1. Experimental Settings

**Implementation Details.** We train SeaS by fine-tuning the pre-trained Stable Diffusion v1-4 [31]. In anomaly image generation experiments, we use 60 normal images and  $\frac{1}{3}$  masked anomaly images for each anomaly type in training. We train one generative model per product, covering all anomaly types. During inference, we generate 1,000 anomaly image-mask pairs for a single anomaly type. More details are given in appendix A.3.

**Datasets.** We conduct experiments on MVTec AD dataset [3], VisA dataset[47], and MVTec 3D AD dataset (only RGB images) [4]. MVTec AD dataset contains 15 product categories, each with up to 8 different anomalies. VisA dataset covers 12 objects in 3 domains. MVTec 3D AD dataset includes 10 product categories, each with up to 4 different anomalies. It contains more challenges, i.e., lighting condition variations, and product pose variations.

**Evaluation Metrics.** For image generation, unlike existing methods [11, 17] that only assess the whole anomaly images, our evaluation contains three levels: anomaly images, normal images, and anomalies, using 4 metrics: (1) Inception Score (IS) and Intra-cluster pairwise LPIPS distance (IC-LPIPS) [29] for authenticity and diversity of anomaly images. (2) KID [5] for authenticity of normal images. (3) IC-LPIPS calculated only in anomaly regions (short for

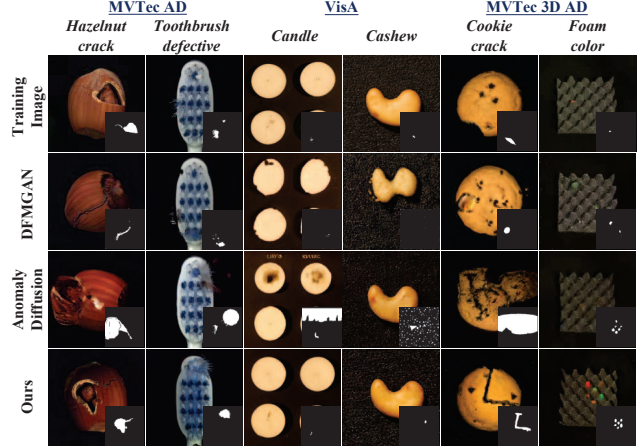


Figure 5. Visualization of the generation results on MVTec AD, VisA and MVTec 3D AD. The sub-image in the lower right corner is the generated mask, none means that the method cannot generate masks.

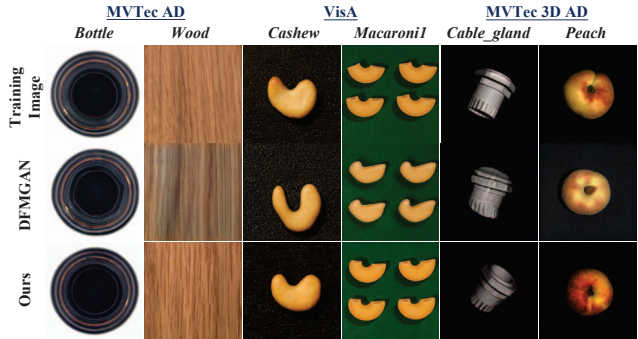


Figure 6. Visualization of the generated normal images on MVTec AD, VisA and MVTec 3D AD.

IC-LPIPS(a)) for the diversity of anomalies. For pixel-level anomaly segmentation and image-level anomaly detection, we use 3 metrics: Area Under Receiver Operator Characteristic curve (AUROC), Average Precision (AP) and  $F_1$ -score at optimal threshold ( $F_1$ -max). We also report Intersection over Union (IoU) for segmentation.

### 4.2. Comparison in Anomaly Image Generation

**Comparison Methods.** For image generation, we compare SeaS with current anomaly image generation methods, like Crop&Paste [23], SDGAN [28], Defect-GAN [42](all without open-source code), DFMGAN [11], and AnomalyDiffusion [17] in terms of fidelity and diversity. For diverse generated anomalies, we combine SeaS-generated anomalies with synthesis-based AD approaches like DRAEM [41] and GLASS [9]. For authentic generated normal images, we use SeaS-generated normal images to augment the training sets of unsupervised AD methods like HVQ-Trans [26], PatchCore [32], and MambaAD [15]. For anomaly image-mask pairs, we generate them with DFMGAN, AnomalyDiffusion, and SeaS, to train segmentation models like BiSeNet V2 [40], UPerNet [38], and LFD [45] respectively. Different from AnomalyDiffusion, which trains one segmentation

Table 1. Comparison on IS and IC-LPIPS on MVTec AD, VisA, and MVTec AD 3D. Bold indicates the best performance.

Methods	MVTec AD				VisA				MVTec 3D AD			
	IS $\uparrow$	IC-LPIPS $\uparrow$	KID $\downarrow$	IC-LPIPS(a) $\uparrow$	IS $\uparrow$	IC-LPIPS $\uparrow$	KID $\downarrow$	IC-LPIPS(a) $\uparrow$	IS $\uparrow$	IC-LPIPS $\uparrow$	KID $\downarrow$	IC-LPIPS(a) $\uparrow$
Crop&Paste[23]	1.51	0.14	-	-	-	-	-	-	-	-	-	-
SDGAN[28]	1.71	0.13	-	-	-	-	-	-	-	-	-	-
Defect-GAN[42]	1.69	0.15	-	-	-	-	-	-	-	-	-	-
DFMGAN[11]	1.72	0.20	0.12	0.14	1.25	0.25	0.24	0.05	1.80	0.29	0.19	0.08
AnomalyDiffusion[17]	1.80	0.32	-	0.12	1.26	0.25	-	0.04	1.61	0.22	-	0.07
<b>SeaS</b>	<b>1.88</b>	<b>0.34</b>	<b>0.04</b>	<b>0.18</b>	<b>1.27</b>	<b>0.26</b>	<b>0.02</b>	<b>0.06</b>	<b>1.95</b>	<b>0.30</b>	<b>0.06</b>	<b>0.09</b>

Table 2. Comparison on combining generated anomalies with synthesis-based anomaly detection methods across multiple datasets.

Segmentation Models	MVTec AD								VisA								MVTec 3D AD							
	Image-level			Pixel-level					Image-level			Pixel-level					Image-level			Pixel-level				
	AUROC	AP	F <sub>1</sub> -max	AUROC	AP	F <sub>1</sub> -max	IoU	AUROC	AP	F <sub>1</sub> -max	AUROC	AP	F <sub>1</sub> -max	IoU	AUROC	AP	F <sub>1</sub> -max	AUROC	AP	F <sub>1</sub> -max	IoU			
DRAEM [41]	98.00	98.45	96.34	97.90	67.89	66.04	<b>60.30</b>	86.28	85.30	81.66	92.92	17.15	22.95	13.57	79.16	90.90	89.78	86.73	14.02	17.00	12.42			
DRAEM + SeaS	<b>98.64</b>	<b>99.40</b>	<b>97.89</b>	<b>98.11</b>	<b>76.55</b>	<b>72.70</b>	58.87	<b>88.12</b>	<b>87.04</b>	<b>83.04</b>	<b>98.45</b>	<b>49.05</b>	<b>48.62</b>	<b>35.00</b>	<b>85.45</b>	<b>93.58</b>	<b>90.85</b>	<b>95.43</b>	<b>20.09</b>	<b>26.10</b>	<b>17.07</b>			
GLASS [9]	99.92	99.98	99.60	99.27	74.09	70.42	57.14	97.68	96.89	93.03	<b>98.47</b>	45.58	48.39	39.92	92.34	96.85	<b>93.37</b>	98.46	48.46	49.13	45.03			
GLASS + SeaS	<b>99.97</b>	<b>99.99</b>	<b>99.81</b>	<b>99.29</b>	<b>76.82</b>	<b>72.38</b>	<b>57.45</b>	<b>97.88</b>	<b>97.39</b>	<b>93.21</b>	98.43	<b>48.06</b>	<b>49.32</b>	<b>40.00</b>	<b>92.95</b>	<b>97.38</b>	<b>93.21</b>	<b>98.73</b>	<b>48.55</b>	<b>49.28</b>	<b>46.02</b>			
Average [9, 41]	98.96	99.22	97.97	98.59	70.99	68.23	<b>58.72</b>	91.98	91.10	87.35	95.70	31.37	35.67	26.75	85.75	93.88	91.58	92.60	31.24	33.07	28.73			
Average(+ SeaS)	<b>99.31</b>	<b>99.70</b>	<b>98.85</b>	<b>98.70</b>	<b>76.69</b>	<b>72.54</b>	58.16	<b>93.00</b>	<b>92.22</b>	<b>88.13</b>	<b>98.44</b>	<b>48.56</b>	<b>48.97</b>	<b>37.50</b>	<b>89.20</b>	<b>95.48</b>	<b>92.03</b>	<b>97.08</b>	<b>34.32</b>	<b>37.69</b>	<b>31.55</b>			

Table 3. Comparison on combining generated normal images with unsupervised anomaly detection methods across multiple datasets.

Segmentation Models	MVTec AD								VisA								MVTec 3D AD							
	Image-level			Pixel-level					Image-level			Pixel-level					Image-level			Pixel-level				
	AUROC	AP	F <sub>1</sub> -max	AUROC	AP	F <sub>1</sub> -max	IoU	AUROC	AP	F <sub>1</sub> -max	AUROC	AP	F <sub>1</sub> -max	IoU	AUROC	AP	F <sub>1</sub> -max	AUROC	AP	F <sub>1</sub> -max	IoU			
HVQ-Trans [26]	96.38	98.09	95.30	<b>97.60</b>	47.95	53.32	<b>45.03</b>	90.11	88.18	84.08	98.10	28.67	35.05	<b>24.03</b>	68.15	84.38	85.20	96.40	17.23	24.59	<b>20.51</b>			
HVQ-Trans + SeaS	<b>97.25</b>	<b>98.48</b>	<b>95.78</b>	97.58	<b>48.53</b>	<b>53.84</b>	44.61	<b>92.12</b>	<b>90.35</b>	<b>86.23</b>	<b>98.15</b>	<b>29.52</b>	<b>36.00</b>	23.60	<b>71.26</b>	<b>90.35</b>	<b>89.23</b>	<b>96.56</b>	<b>19.34</b>	<b>26.40</b>	20.47			
PatchCore [32]	98.63	99.47	98.18	<b>98.37</b>	56.13	58.83	49.45	94.84	95.98	91.69	98.38	48.58	49.69	42.44	83.44	94.89	92.24	98.55	34.52	39.09	39.29			
PatchCore + SeaS	<b>98.64</b>	<b>99.48</b>	<b>98.22</b>	<b>98.37</b>	<b>63.98</b>	<b>64.07</b>	<b>55.43</b>	<b>94.97</b>	<b>96.06</b>	<b>91.81</b>	<b>98.41</b>	<b>48.60</b>	<b>49.72</b>	<b>42.46</b>	<b>83.88</b>	<b>94.97</b>	<b>92.32</b>	<b>98.56</b>	<b>34.65</b>	<b>39.41</b>	<b>39.43</b>			
MambaAD [15]	98.54	99.52	97.77	<b>97.67</b>	56.23	59.34	51.31	94.19	94.44	89.55	98.49	39.27	<b>44.18</b>	<b>37.68</b>	85.92	95.69	92.51	98.57	<b>37.30</b>	<b>41.08</b>	39.44			
MambaAD + SeaS	<b>98.80</b>	<b>99.64</b>	<b>98.40</b>	97.66	<b>56.86</b>	<b>59.70</b>	<b>51.51</b>	<b>94.23</b>	<b>94.65</b>	<b>89.93</b>	<b>98.70</b>	<b>39.33</b>	43.99	36.62	<b>88.67</b>	<b>96.60</b>	<b>93.41</b>	<b>98.74</b>	35.46	<b>39.59</b>	<b>39.51</b>			
Average [15, 26, 32]	97.85	99.03	97.08	<b>97.88</b>	53.44	57.16	48.60	93.05	92.87	88.44	98.32	38.84	42.97	<b>34.72</b>	79.17	91.65	89.98	97.84	29.68	34.92	33.08			
Average(+ SeaS)	<b>98.23</b>	<b>99.20</b>	<b>97.47</b>	97.87	<b>56.46</b>	<b>59.20</b>	<b>50.52</b>	<b>93.77</b>	<b>93.69</b>	<b>89.32</b>	<b>98.42</b>	<b>39.15</b>	<b>43.24</b>	34.23	<b>81.27</b>	<b>93.97</b>	<b>91.65</b>	<b>97.95</b>	<b>29.82</b>	<b>35.13</b>	<b>33.14</b>			

Table 4. Comparison on trained supervised segmentation models for anomaly detection and segmentation across multiple datasets.

Segmentation Models	Generative Models	MVTec AD								VisA								MVTec 3D AD							
		Image-level			Pixel-level					Image-level			Pixel-level					Image-level			Pixel-level				
		AUROC	AP	F <sub>1</sub> -max	AUROC	AP	F <sub>1</sub> -max	IoU	AUROC	AP	F <sub>1</sub> -max	AUROC	AP	F <sub>1</sub> -max	IoU	AUROC	AP	F <sub>1</sub> -max	AUROC	AP	F <sub>1</sub> -max	IoU			
BiSeNet V2 [40]	DFMGAN	90.90	94.43	90.33	94.57	60.42	60.54	45.83	63.07	62.63	66.48	75.91	9.17	15.00	9.66	61.88	81.80	84.44	75.89	15.02	21.73	15.68			
	AnomalyDiffusion	90.08	94.84	91.84	96.27	64.50	62.27	42.89	76.11	77.74	73.13	89.29	34.16	37.93	15.93	61.49	81.35	85.36	<b>92.39</b>	15.15	20.09	14.70			
	SeaS	<b>96.00</b>	<b>98.14</b>	<b>95.43</b>	<b>97.21</b>	<b>69.21</b>	<b>66.37</b>	<b>55.28</b>	<b>85.61</b>	<b>86.64</b>	<b>80.49</b>	<b>96.03</b>	<b>42.80</b>	<b>45.41</b>	<b>25.93</b>	<b>73.60</b>	<b>87.75</b>	<b>85.82</b>	90.41	<b>26.04</b>	<b>32.61</b>	<b>28.55</b>			
UPerNet [38]	DFMGAN	90.74	94.43	90.37	92.33	57.01	56.91	46.64	71.69	71.64	70.70	75.09	12.42	18.52	15.47	67.56	84.53	84.99	75.12	19.54	26.04	18.78			
	AnomalyDiffusion	96.62	98.61	96.21	96.87	69.92	66.95	50.80	83.18	84.08	78.88	95.00	39.92	45.37	20.53	76.56	90.42	87.35	88.48	28.95	35.81	25.04			
	SeaS	<b>98.29</b>	<b>99.20</b>	<b>97.34</b>	<b>97.87</b>	<b>74.42</b>	<b>70.70</b>	<b>61.24</b>	<b>90.34</b>	<b>90.73</b>	<b>84.33</b>	<b>97.01</b>	<b>55.46</b>	<b>55.99</b>	<b>35.91</b>	<b>82.57</b>	<b>92.59</b>	<b>88.72</b>	<b>91.93</b>	<b>38.51</b>	<b>43.53</b>	<b>38.56</b>			
LFD [45]	DFMGAN	91.08	95.40	90.58	94.91	67.06	65.09	45.49	65.38	62.25	66.59	81.21	15.14	18.70	6.44	62.23	82.17	85.38	72.15	9.54	14.29	14.81			
	AnomalyDiffusion	95.15	97.78	94.66	96.30	69.77	66.99	45.77	81.97	82.36	77.35	88.00	30.86	38.56	16.61	77.06	89.44	87.20	<b>92.68</b>	24.29	32.74	19.90			
	SeaS	<b>95.88</b>	<b>97.89</b>	<b>95.15</b>	<b>98.09</b>	<b>77.15</b>	<b>72.52</b>	<b>56.47</b>	<b>83.07</b>	<b>82.88</b>	<b>77.24</b>	<b>92.91</b>	<b>43.87</b>	<b>46.46</b>	<b>26.37</b>	<b>78.96</b>	<b>91.22</b>	<b>87.28</b>	91.61	<b>40.25</b>	<b>43.47</b>	<b>39.00</b>			
Average	DFMGAN	90.91	94.75	90.43	93.94	61.50	60.85	45.99	66.71	65.51	67.92	77.40	12.24	17.41	10.52	63.89	83.83	84.94	74.39	14.70	20.69	16.42			
	AnomalyDiffusion	93.95	97.08	94.24	96.48	68.06	65.40	46.49	80.42	81.39	76.45	90.76	34.98	40.62	17.69	71.70	87.07	86.64	91.18	22.80	29.55	19.88			
	SeaS	<b>96.72</b>	<b>98.41</b>	<b>95.97</b>	<b>97.72</b>	<b>73.59</b>	<b>69.86</b>	<b>57.66</b>	<b>86.34</b>	<b>86.75</b>	<b>80.69</b>	<b>95.32</b>	<b>47.38</b>	<b>49.29</b>	<b>29.40</b>	<b>78.38</b>	<b>90.52</b>	<b>87.27</b>	<b>91.32</b>	<b>34.93</b>	<b>39.87</b>	<b>35.37</b>			

model per product, we train a unified supervised segmentation model for all products, which is more challenging.

**Anomaly image generation quality.** In Tab. 1, we compare SeaS with some state-of-the-art anomaly image generation methods on fidelity (IS and KID) and diversity (IC-LPIPS and IC-LPIPS(a)). SeaS outperforms other methods in IS and IC-LPIPS, showing superior fidelity and diversity. It also excels in generating authentic normal images and diverse anomalies. Compared to AnomalyDiffusion, which cannot generate normal images, SeaS leads in IC-LPIPS(a). SeaS also surpasses DFMGAN in both KID and IC-LPIPS(a). We exhibit the generated anomaly images in Fig. 5, SeaS-generated anomaly images have higher

fidelity (e.g., *hazelnut\_crack*). Compared with other methods, SeaS can generate images with different types, colors, and shapes of anomalies rather than overfitting to the training images (e.g., *foam\_color*). SeaS-generated masks are also precisely aligned with the anomaly regions (e.g., *tooth-brush\_defective*). We also present the authentic generated normal images in Fig. 6. More qualitative and quantitative anomaly image generation results are in appendix A.6.

**Combining generated anomalies with synthesis-based AD methods.** We replace the synthesized pseudo-anomalies in DRAEM [41] and GLASS [9] with SeaS-generated anomalies. As shown in Tab. 2, SeaS-generated anomalies, which offer sufficient diversity, consistently im-

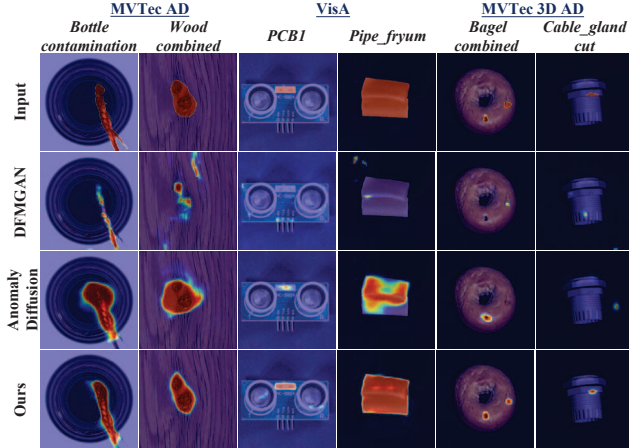


Figure 7. Qualitative supervised anomaly segmentation results with BiSeNet V2 on MVTec AD.

prove synthesis-based AD methods by suppressing false negatives, leading to better performance across multiple datasets. More training details are given in appendix A.3.

#### Combining generated normal images with AD methods.

We use SeaS-generated normal images to supplement the training sets of existing state-of-the-art unsupervised AD methods, the results are given in Tab. 3. Using SeaS-generated normal images with minor local variations and global consistency, unsupervised AD methods reduce false positives and perform well across multiple metrics, improving industrial anomaly detection on various datasets. More training details are given in appendix A.3.

#### Training supervised segmentation models for anomaly segmentation and detection.

We generate 1,000 image-mask pairs for each anomaly type and use them, along with all normal images in the original training sets, to train a unified supervised segmentation model. The models are tested on the remaining images not included in the training set. All methods are trained using the same number of images and the training settings, detailed in appendix A.4. As shown in Tab. 4, the segmentation results consistently demonstrate that our method outperforms others across all the segmentation models, with average IoU improvements of 11.17% (MVTec AD), 11.71% (VisA), and 15.49% (MVTec 3D AD). Segmentation anomaly maps are shown in Fig. 7. Using our generated image-mask pairs to train BiSeNet V2, there are fewer false positives in *wood\_combined* and fewer false negatives in *bottle\_contamination*. We also use the maximum value of the segmentation anomaly map as the image-level anomaly score for anomaly detection, achieving gains of 2.77% (MVTec AD), 5.92% (VisA), and 6.68% (MVTec 3D AD) in image-AUROC. More qualitative comparison results are in appendix A.7 and appendix A.8.

#### 4.3. Ablation Study

We train additional models to assess the effect of each component: (a) the model with predefined typical text prompt with fixed generic semantic words (short for with TP in Tab.

Table 5. Ablation on the generation model.

Method	Metrics					
	IS	IC-L	AUROC	AP	$F_1$ -max	IoU
(a) with TP	1.72	0.33	94.72	57.16	55.67	50.46
(b) w/o Mixed	1.79	0.32	95.82	66.07	64.50	53.11
(c) w/o NA	1.67	0.31	96.20	66.03	64.09	53.97
(d) w/o ST	1.86	0.33	96.44	67.73	65.23	54.99
(e) All (Ours)	<b>1.88</b>	<b>0.34</b>	<b>97.21</b>	<b>69.21</b>	<b>66.37</b>	<b>55.28</b>

5); (b) the model without mixing different types of anomaly images in the same product; (c) the model without NA loss; (d) the model without the second term of DA loss in Eq. 3 (short for ST in Tab. 5); (e) our complete model. We use these models to generate 1,000 anomaly image-mask pairs per anomaly type and train BiSeNet V2 for supervised anomaly segmentation. In Tab. 5, the results show that omitting any component leads to a decrease in fidelity and diversity of the generated images, as well as in the segmentation results. These validate the effectiveness of the components we proposed. More ablation studies on SeaS are shown in appendix A.5.

**Refined Mask Prediction branch.** To verify the validity of the components in the RMP branch, we conduct ablation studies on MRM, the progressive manner to refine coarse feature (short for PM in Tab. 6) and coarse mask supervision (short for CMS in Tab. 6). 1) the model only with CMS, which means we do not use MRM to fuse the high-resolution features in RMP, but directly obtain the mask from the coarse features  $\hat{F} \in \mathbb{R}^{64 \times 64 \times 192}$  through convolution and bilinear interpolation upsampling; 2) the model with MRM; 3) the model utilizing three MRMs in a progressive manner to refine coarse features; 4) our complete model. We report the BiSeNet V2 results in Tab. 6, which demonstrates that each component in the RMP is indispensable for downstream supervised anomaly segmentation. More ablation studies about RMP are in appendix A.5.

Table 6. Ablation on the RMP branch.

Method			Metrics			
MRM	PM	CMS	AUROC	AP	$F_1$ -max	IoU
		✓	97.00	65.28	62.56	53.93
✓			94.54	60.52	59.06	49.42
✓	✓		94.04	62.04	59.82	50.44
✓	✓	✓	<b>97.21</b>	<b>69.21</b>	<b>66.37</b>	<b>55.28</b>

## 5. Conclusion

In this paper, we propose a unified generation method named SeaS. We explore an implicit characteristic that anomalies exhibit high variability, while normal products maintain global consistency. We design a Separation and Sharing Fine-tuning strategy to model different variations of normal products and anomalies, enabling the Refined Mask Prediction branch to predict accurate masks with discriminative features. Our method greatly improves synthesis-based and supervised AD methods, and empowers supervised segmentation models.

## 6. Acknowledgement

This work was supported by the National Natural Science Foundation of China under Grant No.62176098. The computation is completed in the HPC Platform of Huazhong University of Science and Technology.

## References

- [1] Omri Avrahami, Kfir Aberman, Ohad Fried, Daniel Cohen-Or, and Dani Lischinski. Break-a-scene: Extracting multiple concepts from a single image. In *SIGGRAPH Asia 2023 Conference Papers*, pages 1–12, 2023. [2](#)
- [2] Shane Barratt and Rishi Sharma. A note on the inception score. *arXiv preprint arXiv:1801.01973*, 2018. [12](#)
- [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. [6](#)
- [4] Paul Bergmann., Xin Jin., David Sattlegger., and Carsten Steger. The mvtec 3d-ad dataset for unsupervised 3d anomaly detection and localization. In *Proceedings of the 17th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications*, pages 202–213, 2022. [6](#)
- [5] Mikołaj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. In *International Conference on Learning Representations*, 2018. [6](#), [12](#)
- [6] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023. [2](#)
- [7] Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. *ACM Transactions on Graphics*, 42(4):1–10, 2023. [3](#)
- [8] Hong Chen, Yipeng Zhang, Xin Wang, Xuguang Duan, Yuwei Zhou, and Wenwu Zhu. Disenbooth: Disentangled parameter-efficient tuning for subject-driven text-to-image generation. In *International Conference on Learning Representations*, 2024. [2](#)
- [9] Qiyu Chen, Huiyuan Luo, Chengkan Lv, and Zhengtao Zhang. A unified anomaly synthesis strategy with gradient ascent for industrial anomaly detection and localization. In *European Conference on Computer Vision*, pages 37–54. Springer, 2024. [1](#), [6](#), [7](#), [12](#)
- [10] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017. [2](#)
- [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. [1](#), [2](#), [6](#), [7](#), [12](#), [13](#), [17](#), [26](#)
- [12] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *International Conference on Learning Representations*, 2022. [2](#), [35](#)
- [13] Guan Gui, Bin-Bin Gao, Jun Liu, Chengjie Wang, and Yunsheng Wu. Few-shot anomaly-driven generation for anomaly classification and segmentation. In *European Conference on Computer Vision*, pages 210–226, 2024. [2](#)
- [14] Ligong Han, Yinxiao Li, Han Zhang, Peyman Milanfar, Dimitris Metaxas, and Feng Yang. Svdif: Compact parameter space for diffusion fine-tuning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7323–7334, 2023. [2](#)
- [15] Haoyang He, Yuhu Bai, Jiangning Zhang, Qingdong He, Hongxu Chen, Zhenye Gan, Chengjie Wang, Xiangtai Li, Guanzhong Tian, and Lei Xie. Mambaad: Exploring state space models for multi-class unsupervised anomaly detection. *Advances in Neural Information Processing Systems*, 37:71162–71187, 2025. [1](#), [6](#), [7](#), [12](#)
- [16] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or. Prompt-to-prompt image editing with cross-attention control. In *International Conference on Learning Representations*, 2022. [3](#)
- [17] Teng Hu, Jiangning Zhang, Ran Yi, Yuzhen Du, Xu Chen, Liang Liu, Yabiao Wang, and Chengjie Wang. Anomalydiffusion: Few-shot anomaly image generation with diffusion model. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2024. [1](#), [2](#), [6](#), [7](#), [12](#), [13](#), [17](#), [26](#), [35](#)
- [18] Jongheon Jeong, Yang Zou, Taewan Kim, Dongqing Zhang, Avinash Ravichandran, and Onkar Dabeer. Winclip: Zero/few-shot anomaly classification and segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19606–19616, 2023. [3](#)
- [19] Chen Jin, Ryutaro Tanno, Amrutha Saseendran, Tom Diefte, and Philip Teare. An image is worth multiple words: Learning object level concepts using multi-concept prompt learning. In *International Conference on Machine Learning*, 2024. [2](#)
- [20] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1931–1941, 2023. [2](#)
- [21] 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](#)
- [22] Daiqing Li, Huan Ling, Seung Wook Kim, Karsten Kreis, Sanja Fidler, and Antonio Torralba. Bigdatasetgan: Synthesizing imagenet with pixel-wise annotations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21330–21340, 2022. [2](#)
- [23] Dongyun Lin, Yanpeng Cao, Wenbin Zhu, and Yiqun Li. Few-shot defect segmentation leveraging abundant defect-free training samples through normal background regularization and crop-and-paste operation. In *2021 IEEE International Conference on Multimedia and Expo*, pages 1–6, 2021. [6](#), [7](#), [17](#)

- [24] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2980–2988, 2017. 5
- [25] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018. 12
- [26] Ruiying Lu, YuJie Wu, Long Tian, Dongsheng Wang, Bo Chen, Xiyang Liu, and Ruimin Hu. Hierarchical vector quantized transformer for multi-class unsupervised anomaly detection. *Advances in Neural Information Processing Systems*, 36:8487–8500, 2023. 1, 6, 7, 12
- [27] Quang Nguyen, Truong Vu, Anh Tran, and Khoi Nguyen. Dataset diffusion: Diffusion-based synthetic data generation for pixel-level semantic segmentation. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- [28] Shuanlong Niu, Bin Li, Xinggong Wang, and Hui Lin. Defect image sample generation with gan for improving defect recognition. *IEEE Transactions on Automation Science and Engineering*, 17(3):1611–1622, 2020. 2, 6, 7, 17
- [29] Utkarsh Ojha, Yijun Li, Jingwan Lu, Alexei A. Efros, Yong Jae Lee, Eli Shechtman, and Richard Zhang. Few-shot image generation via cross-domain correspondence. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10743–10752, 2021. 6, 12
- [30] Patrick Pérez, Michel Gangnet, and Andrew Blake. Poisson image editing. In *ACM SIGGRAPH 2003*, pages 313–318, 2003. 2
- [31] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 2, 3, 6
- [32] 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. 1, 6, 7, 12
- [33] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22500–22510, 2023. 2
- [34] Hannah M Schlüter, Jeremy Tan, Benjamin Hou, and Bernhard Kainz. Natural synthetic anomalies for self-supervised anomaly detection and localization. In *European Conference on Computer Vision*, pages 474–489. Springer, 2022. 2
- [35] Weijia Wu, Yuzhong Zhao, Hao Chen, Yuchao Gu, Rui Zhao, Yefei He, Hong Zhou, Mike Zheng Shou, and Chunhua Shen. Datasetdm: Synthesizing data with perception annotations using diffusion models. *Advances in Neural Information Processing Systems*, 36:54683–54695, 2023. 3
- [36] Weijia Wu, Yuzhong Zhao, Mike Zheng Shou, Hong Zhou, and Chunhua Shen. Diffumask: Synthesizing images with pixel-level annotations for semantic segmentation using diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1206–1217, 2023. 3
- [37] Guangxuan Xiao, Tianwei Yin, William T Freeman, Frédo Durand, and Song Han. Fastcomposer: Tuning-free multi-subject image generation with localized attention. *arXiv preprint arXiv:2305.10431*, 2023. 2
- [38] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *Proceedings of the European Conference on Computer Vision*, pages 418–434, 2018. 6, 7, 12, 26
- [39] Jinheng Xie, Yuexiang Li, Yawen Huang, Haozhe Liu, Wentian Zhang, Yefeng Zheng, and Mike Zheng Shou. Boxdiff: Text-to-image synthesis with training-free box-constrained diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7452–7461, 2023. 3
- [40] Changqian Yu, Changxin Gao, Jingbo Wang, Gang Yu, Chunhua Shen, and Nong Sang. Bisenet v2: Bilateral network with guided aggregation for real-time semantic segmentation. *International Journal of Computer Vision*, 129:3051–3068, 2021. 6, 7, 12, 26
- [41] Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. 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, 6, 7, 12, 35
- [42] Gongjie Zhang, Kaiwen Cui, Tzu-Yi Hung, and Shijian Lu. Defect-gan: High-fidelity defect synthesis for automated defect inspection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2524–2534, 2021. 2, 6, 7, 17
- [43] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 2
- [44] Yuxuan Zhang, Huan Ling, Jun Gao, Kangxue Yin, Jean-Francois Lafleche, Adela Barriuso, Antonio Torralba, and Sanja Fidler. Datasetgan: Efficient labeled data factory with minimal human effort. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10145–10155, 2021. 2
- [45] Huan Zhou, Feng Xue, Yucong Li, Shi Gong, Yiqun Li, and Yu Zhou. Exploiting low-level representations for ultra-fast road segmentation. *IEEE Transactions on Intelligent Transportation Systems*, 2024. 6, 7, 12, 26
- [46] Qihang Zhou, Guansong Pang, Yu Tian, Shibo He, and Jiming Chen. Anomalyclip: Object-agnostic prompt learning for zero-shot anomaly detection. In *International Conference on Learning Representations*, 2024. 3
- [47] Yang Zou, Jongheon Jeong, Latha Pemula, Dongqing Zhang, and Onkar Dabeer. Spot-the-difference self-supervised pre-training for anomaly detection and segmentation. In *European Conference on Computer Vision*, pages 392–408. Springer, 2022. 6