

# SEAS: FEW-SHOT INDUSTRIAL ANOMALY IMAGE GENERATION WITH SEPARATION AND SHARING FINE-TUNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Current segmentation methods typically require many training images and precise masks, while insufficient anomaly images hinder their application in industrial scenarios. To address such an issue, we explore producing diverse anomalies and accurate pixel-wise annotations. By observing the real production lines, we find that anomalies vary randomly in shape and appearance, whereas products hold globally consistent patterns with slight local variations. Such a characteristic inspires us to develop a Separation and Sharing Fine-tuning (SeaS) approach using only a few abnormal and some normal images. Firstly, we propose the Unbalanced Abnormal (UA) Text Prompt tailored to industrial anomaly generation, consisting of one product token and several anomaly tokens. Then, for anomaly images, we propose a Decoupled Anomaly Alignment (DA) loss to bind the attributes of the anomalies to different anomaly tokens. Re-blending such attributes may produce never-seen anomalies, achieving a high diversity of anomalies. For normal images, we propose a Normal-image Alignment (NA) loss to learn the products' key features that are used to synthesize products with both global consistency and local variations. The two training processes are separated but conducted on a shared U-Net. Finally, SeaS produces high-fidelity annotations for the generated anomalies by fusing discriminative features of U-Net and high-resolution VAE features. The extensive evaluations on the challenging MVTec AD and MVTec 3D AD dataset (RGB images) demonstrate the effectiveness of our approach. For anomaly image generation, on MVTec AD dataset, we achieve 1.88 on IS and 0.34 on IC-LPIPS, while on the MVTec 3D AD dataset, we obtain 1.95 on IS and 0.30 on IC-LPIPS. For the downstream task, by using our generated anomaly image-mask pairs, three common segmentation methods achieve an average 11.17% improvement on IoU on MVTec AD dataset, and a 15.49% enhancement in IoU on the MVTec 3D AD dataset. The source code will be released publicly available.

## 1 INTRODUCTION

Existing segmentation approaches require a large number of anomaly images with mask annotations, while the scarcity of anomaly images obstructs their application in industrial scenarios. To solve this problem, generative methods for industrial scenarios have emerged to expand the training set of segmentation models.

To the best of our knowledge, generation approaches (Zavrtanik et al., 2021) in industrial scenarios can be broadly classified into two categories: **Anomaly Generation (AG)** and **Anomaly Image Generation (AIG)**. AG methods (Li et al., 2021; Zavrtanik et al., 2021; Schlüter et al., 2022; Hu et al., 2024) generate anomalies only and merge them into the real normal images using different strategies, e.g., CutPaste (Li et al., 2021) pastes a cropped normal region to normal images, which simulates anomalies by misalignment. AnomalyDiffusion (Hu et al., 2024) generates anomalies by a diffusion model, and edits anomalies onto the normal images guided by the anomaly masks, as shown in Fig. 1(a). However, AG methods require anomaly masks as inputs, which easily suffer from low fidelity and consistency in generation if these masks are unreasonably positioned. In contrast to AG, as shown in Fig. 1(b), AIG approaches take a step further, generating anomalies and the industrial products that they lie in simultaneously. Therefore, AIG faces greater challenges

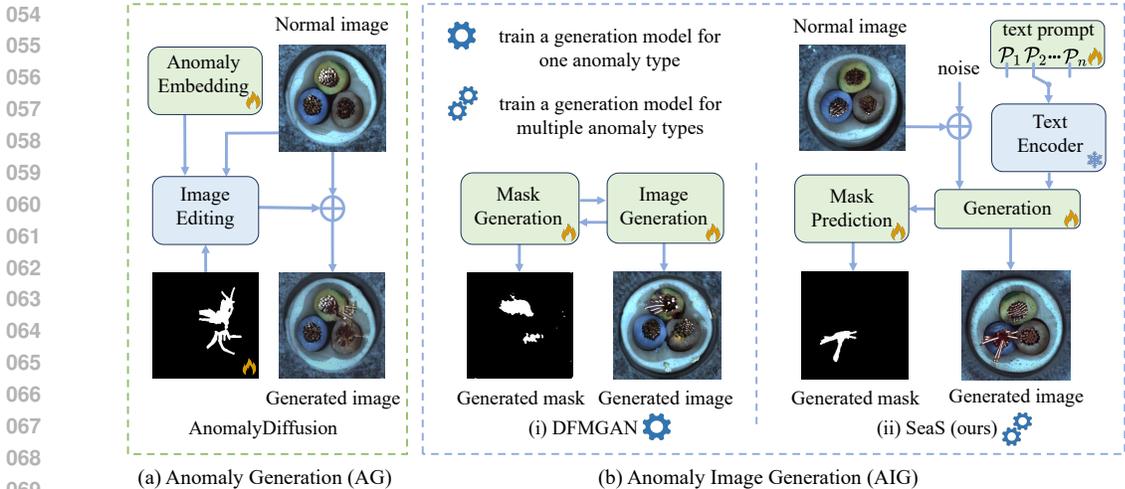


Figure 1: (a) **Anomaly Generation (AG)** method only generates anomaly and edits it onto normal images guided by the anomaly masks. (b) **Anomaly Image Generation (AIG)** methods generate anomalies and the industrial products that they lie in simultaneously. (i) DFMGAN trains a generation model for each anomaly type. (ii) SeaS trains a shared generation model for multiple anomaly types.

due to the requirement of high-fidelity of both anomalies and products. Previously, AIG methods (Duan et al., 2023; Zhang et al., 2021a; Niu et al., 2020) are mainly based on Generative Adversarial Network (GAN). However, they are limited by insufficient generalization of GAN, lacking specific control over products/anomalies, and inaccurate masks.

In real industrial manufacturing, the products in an individual production line are almost similar to each other, while the anomalies are unforeseeable in shape and appearance. **Such an observation reveals a differentiated characteristic, i.e., the products satisfy global consistency with minor variations in local details, while the anomalies hold randomness**, which is rarely discussed in existing AG or AIG approaches. Motivated by such a characteristic, we propose a **Separation and Sharing Fine-tuning** method, short by SeaS, a controllable AIG method based on Stable Diffusion (Rombach et al., 2022). The key idea is to employ Unbalanced Abnormal (UA) Text Prompts containing a set of tokens that characterize products and anomalies separately, so that the anomaly tokens align with the anomaly semantics for diverse generations, and a product token expresses a globally consistent product surface. Specifically, to learn highly-diverse anomalies, we first propose a Decoupled Anomaly Alignment (DA) loss to bind the attributes of the anomalies to different anomaly tokens. Recombining the decoupled attributes may produce anomalies that have never been seen in the training dataset, therefore increasing the diversity of the generated anomalies. Secondly, to learn globally-consistent patterns from products, we propose the Normal-image Alignment (NA) loss. It enables the network to learn the key features of the product from normal images and fine-tune a learnable embedding. Such an embedding ensures the preservation of global consistency amidst local detail variations. Thirdly, **according to the experimental analysis, we find that existing methods leverage the low-resolution features to predict the mask, which may introduce a large amount of boundary uncertainty**. Thus, we propose a Refined Mask Prediction (RMP) branch to produce pixel-wise anomaly annotations for other downstream tasks. It combines the discriminative U-Net features and high-resolution VAE features to generate accurate and crisp masks in a progressive way. Extensive experiments on AIG and downstream anomaly segmentation tasks show that SeaS outperforms the existing industrial anomaly generation methods. On MVTEC AD dataset, our model achieves **1.88** on IS metric and **0.34** on IC-LPIPS. Furthermore, training on the image-mask pairs generated by SeaS, the downstream segmentation models achieve improvements of average **+5.53%** AP and **+11.17%** IoU. On MVTEC AD 3D dataset (RGB images), our method attains **1.95** on IS metric and **0.30** on IC-LPIPS. Using the image-mask pairs generated by SeaS to train the downstream segmentation models, we exhibit average improvements of **+12.13%** AP and **+15.49%** IoU.

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

- We reveal different characteristics of products and anomalies, which motivates us to propose SeaS, a novel AIG method. It independently learns products and anomalies on a shared U-Net and ensures the randomness of anomalies and global consistency of products.
- We propose a Refined Mask Prediction branch to produce accurate and crisp pixel-wise annotations for generated anomalies, which combines the advantages of the discriminative U-Net features and the high-resolution VAE decoder features.
- Extensive experiments on anomaly image generation and downstream anomaly segmentation tasks show that SeaS outperforms the existing industrial anomaly generation methods.

## 2 RELATED WORK

**Anomaly Image Generation.** Early non-generative methods (DeVries & Taylor, 2017; Li et al., 2021; Zavrtnik et al., 2021) use data augmentation techniques to create anomaly images. Data augmentation techniques lack consistency in anomaly images, resulting in low fidelity. AG methods (Li et al., 2021; Zavrtnik et al., 2021; Schlüter et al., 2022; Hu et al., 2024; Gui et al., 2024) only generate anomalies and merge them into the real normal images. NSA (Schlüter et al., 2022) uses Poisson Image Editing (Pérez et al., 2003) to facilitate the fusion of the cropped normal region. However, AG methods require anomaly masks as inputs. If these masks are positioned in an unreasonable manner, the generated images will have low fidelity and consistency. Previous AIG methods (Duan et al., 2023; Zhang et al., 2021a; Niu et al., 2020) are mainly based on GAN. Defect-GAN (Zhang et al., 2021a) cannot generate masks. The masks produced by DFMGAN (Duan et al., 2023) often do not align accurately with anomalies, limiting their utility in training segmentation models. We propose a controllable AIG model based on Stable Diffusion to generate high-fidelity and diverse anomaly images with accurate masks.

**Fine-tuning Diffusion Models.** Fine-tuning is a potent strategy for enhancing specific capabilities of pre-trained diffusion models (Gal et al., 2022; Zhang et al., 2023b; Brooks et al., 2023). Personalized methods (Ruiz et al., 2023; Gal et al., 2022; Chen et al., 2024) 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 (Kumari et al., 2023; Xiao et al., 2023; Avrahami et al., 2023; Han et al., 2023; Jin et al., 2024) use cross-attention maps to align embeddings with individual concepts in the image. Nevertheless, they do not consider the diversity requirements between different concepts, which is important for industrial anomaly image generation. Thus, we propose a separation and sharing fine-tuning strategy for the different diversity needs of anomalies and 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 GANs (Zhang et al., 2021b; Li et al., 2022). However, these approaches do not guarantee the generation of accurate masks for exceedingly small datasets. Based on Stable Diffusion (Rombach et al., 2022), some recent methods, i.e., DiffuMask (Wu et al., 2023b), DatasetDM (Wu et al., 2023a) and DatasetDiffusion (Nguyen et al., 2024), 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 the discriminative features of U-Net decoder 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 products and anomalies separately. Subsequently, we propose the Decoupled Anomaly Alignment (DA) loss to bind the anomaly image regions to the anomaly tokens, and leverage Normal-image Alignment (NA) loss to empower the product token to express 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

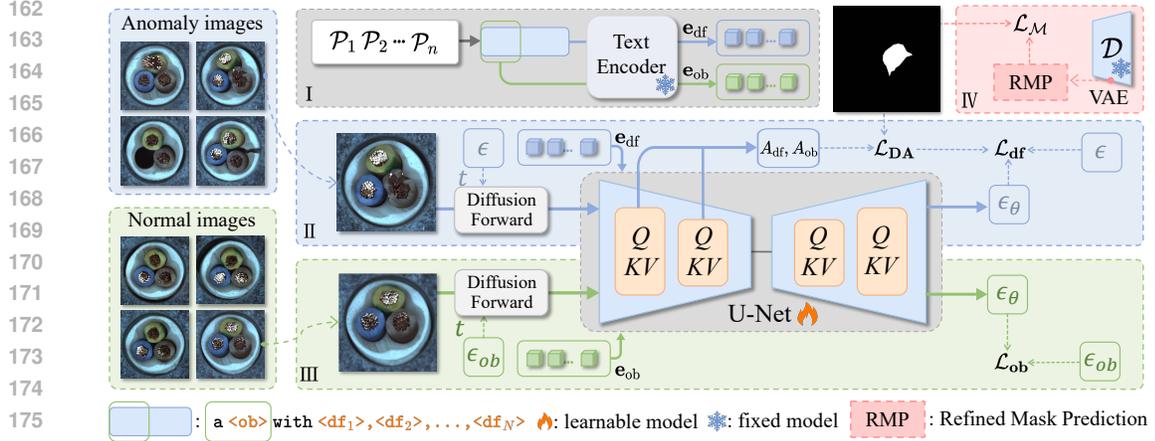


Figure 2: **Overall framework of SeaS.** It consists of four parts: (I) the Unbalanced Abnormal Text Prompt, (II) the Decoupled Anomaly Alignment for aligning the anomaly tokens  $\langle df_n \rangle$  to the anomaly area of abnormal images, (III) the Normal-image Alignment for maintaining authenticity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

branch to generate accurate masks corresponding to the generated anomaly images in Sec. 3.3. Finally, we detail the generation of the abnormal image-mask pairs in Sec. 3.4.

### 3.1 PRELIMINARIES

**Stable Diffusion.** Given an input image  $x_0$ , Stable Diffusion (Rombach et al., 2022) 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 (Hertz et al., 2022; Chefer et al., 2023; Xie et al., 2023). 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 (Jeong et al., 2023), or damaged bottle (Zhou et al., 2024b), is suboptimal for industrial anomaly generation. The fixed generic semantic words, e.g., damaged, defect, may fail to align with a few training images that contain specific defect types. Therefore, we design the Unbalanced Abnormal (UA) Text Prompt for each anomaly type of each product, i.e.,

$$\mathcal{P} = \text{a } \langle ob \rangle \text{ with } \langle df_1 \rangle, \langle df_2 \rangle, \dots, \langle df_N \rangle$$

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

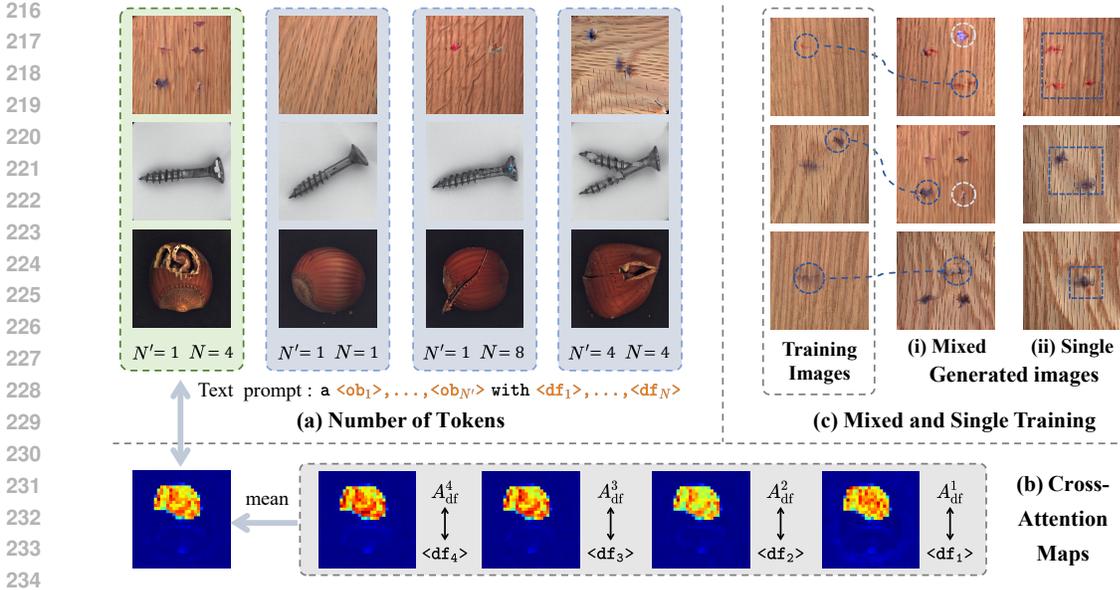


Figure 3: (a) Generated images with the different number of tokens. (b) Cross-attention maps. (c) Examples of diverse generated images.

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 the normal product, while multiple  $\langle df_n \rangle$  are necessary for controlling the generation of the 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 product, which may reduce the authenticity of the 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.,

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

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 number 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 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  $\mathbf{e}_{df} \in \mathbb{R}^{Z \times C_1}$  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 may reduce the authenticity of the generated 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 product 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 images, the objective of Separation and Sharing Fine-tuning is expressed 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 and some normal images with their corresponding masks. 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 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 is 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. 23 in appendix A.8.

### 3.3 REFINED MASK PREDICTION

High-fidelity pixel-wise annotations of anomalies play an important role in boosting segmentation models. However, existing methods, such as DFMGAN (Duan et al., 2023) and AnomalyDiffusion (Hu et al., 2024), produce anomaly masks that are not tightly matched with generated anomalies, which is insufficient for training segmentation model. To address this issue, we design a cascaded Refined Mask Prediction (RMP) branch, which is grafted onto the U-Net trained according to 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 it 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 features 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

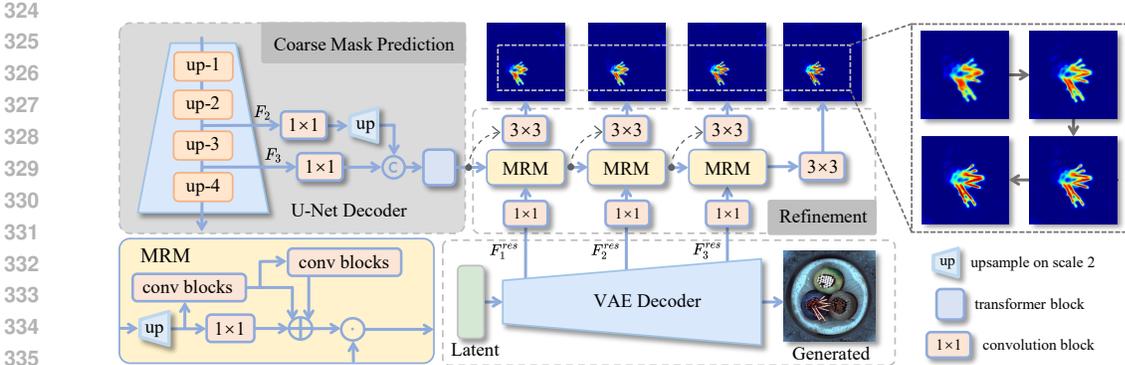


Figure 4: The **Refined Mask Prediction (RMP)** branch during inference. The Coarse Feature Extraction utilizes features from the up-2 and up-3 layers of U-Net Decoder to extract coarse features. The cascaded Mask Refinement Module (MRM) further obtains the mask accurately aligned with the anomaly with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

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 output 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 (Lin et al., 2017).  $\mathbf{M}_{ob} \in \mathbb{R}^{64 \times 64 \times 1}$  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. 27 and Fig. 16 in appendix A.8.

### 3.4 INFERENCE

During the generation of the abnormal image-mask pairs, aiming further to ensure the global consistency of the abnormal image, 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  $z_0$ . Next,  $z_0$  is input into the U-Net for noise prediction, with the process guided by the conditioning vector  $\mathbf{e}_{df}$  (mentioned in Eq. 4). 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 segmentation models is shown in Tab. 29 in appendix A.8. In the last denoising step, the output of the generation model is used as the generated abnormal image.

Table 1: Comparison on IS and IC-LPIPS on MVTec AD. Bold indicates the best performance, while underlined denotes the second-best result.

Category	CDC		Crop& Paste		SDGAN		Defect-GAN		DFMGAN		Anomaly Diffusion		Ours	
	IS $\uparrow$	IC-L $\uparrow$	IS $\uparrow$	IC-L $\uparrow$	IS $\uparrow$	IC-L $\uparrow$								
bottle	1.52	0.04	1.43	0.04	1.57	0.06	1.39	0.07	<u>1.62</u>	0.12	1.58	<u>0.19</u>	<b>1.78</b>	<b>0.21</b>
cable	1.97	0.19	1.74	0.25	1.89	0.19	1.70	0.22	1.96	0.25	<b>2.13</b>	<u>0.41</u>	<u>2.09</u>	<b>0.42</b>
capsule	1.37	0.06	1.23	0.05	1.49	0.03	<b>1.59</b>	0.04	<b>1.59</b>	0.11	<b>1.59</b>	<u>0.21</u>	<u>1.56</u>	<b>0.26</b>
carpet	<b>1.25</b>	0.03	1.17	0.11	1.18	0.11	<u>1.24</u>	0.12	1.23	0.13	1.16	<u>0.24</u>	1.13	<b>0.25</b>
grid	1.97	0.07	2.00	0.12	1.95	0.10	2.01	0.12	1.97	<u>0.13</u>	<u>2.04</u>	<b>0.44</b>	<b>2.43</b>	<b>0.44</b>
hazelnut	<u>1.97</u>	0.05	1.74	0.21	1.85	0.16	1.87	0.19	1.93	<u>0.24</u>	<b>2.13</b>	<b>0.31</b>	1.87	<b>0.31</b>
leather	1.80	0.07	1.47	0.14	2.04	0.12	<b>2.12</b>	0.14	<u>2.06</u>	0.17	1.94	<b>0.41</b>	2.03	0.40
metal_nut	1.55	0.04	1.56	0.15	1.45	0.28	1.47	0.30	1.49	<b>0.32</b>	<b>1.96</b>	0.30	<u>1.64</u>	<u>0.31</u>
pill	1.56	0.06	1.49	0.11	1.61	0.07	1.61	0.10	<b>1.63</b>	0.16	1.61	<u>0.26</u>	<u>1.62</u>	<b>0.33</b>
screw	1.13	0.11	1.12	0.16	1.17	0.10	1.19	0.12	1.12	0.14	<u>1.28</u>	<u>0.30</u>	<b>1.52</b>	<b>0.31</b>
tile	2.10	0.12	1.83	0.20	2.53	0.21	2.35	0.22	2.39	0.22	<u>2.54</u>	<b>0.55</b>	<b>2.60</b>	<u>0.50</u>
toothbrush	1.63	0.06	1.30	0.08	1.78	0.03	<u>1.85</u>	0.03	1.82	0.18	1.68	<u>0.21</u>	<b>1.96</b>	<b>0.25</b>
transistor	1.61	0.13	1.39	0.15	<b>1.76</b>	0.13	1.47	0.13	<u>1.64</u>	<u>0.25</u>	1.57	<b>0.34</b>	1.51	<b>0.34</b>
wood	2.05	0.03	1.95	0.23	2.12	0.25	2.19	0.29	2.12	0.35	<u>2.33</u>	<u>0.37</u>	<b>2.77</b>	<b>0.46</b>
zipper	1.30	0.05	1.23	0.11	1.25	0.10	1.25	0.10	1.29	<u>0.27</u>	<u>1.39</u>	0.25	<b>1.63</b>	<b>0.30</b>
Average	1.65	0.07	1.51	0.14	1.71	0.13	1.69	0.15	1.72	0.20	<u>1.80</u>	<u>0.32</u>	<b>1.88</b>	<b>0.34</b>

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETTINGS

**Implementation Details.** We train SeaS by fine-tuning the pre-trained Stable Diffusion v1-4 (Rom-bach et al., 2022). In AIG experiments, we use 60 normal images and 1/3 abnormal images with their corresponding masks per anomaly type for training. During inference, we generate 1,000 abnormal image-mask pairs for a single anomaly type. More details are given in appendix A.3.

**Datasets.** We conduct experiments on MVTec AD dataset (Bergmann et al., 2019) and MVTec 3D AD dataset (only RGB images) (Bergmann et al., 2022). MVTec AD dataset contains 15 product categories, each with up to 8 different anomalies, making it suitable for simulating real-world industrial scenarios. MVTec 3D AD dataset includes 10 product categories, each with up to 4 different anomalies. It contains more challenges, i.e., lighting condition variations, product pose variations. Due to the page limitation, results on MVTec 3D AD dataset are given in appendix A.4.

**Evaluation Metrics.** For AIG, we leverage 2 metrics: the Inception Score (IS) and the Intra-cluster pairwise LPIPS distance (IC-LPIPS) (Ojha et al., 2021). The scarcity of abnormal images hampers the reliability of FID (Heusel et al., 2017) and KID (Bińkowski et al., 2018), as overfitted model (Duan et al., 2023) achieves higher scores. For pixel-level anomaly segmentation and image-level anomaly detection, we report 3 metrics: Area Under Receiver Operator Characteristic curve (AU-ROC), Average Precision (AP) and  $F_1$ -score at optimal threshold ( $F_1$ -max). In addition, we report Intersection over Union (IoU) for segmentation.

### 4.2 COMPARISON IN ANOMALY IMAGE GENERATION

**Comparison Methods.** We compare SeaS with the current AG and AIG methods on generation fidelity and diversity, such as CDC (Ojha et al., 2021), Crop&Paste (Lin et al., 2021), SDGAN (Niu et al., 2020), Defect-GAN (Zhang et al., 2021a), DFMGAN (Duan et al., 2023) and AnomalyDiffusion (Hu et al., 2024). Then we use Crop&Paste, DRAEM (Zavrtanik et al., 2021), DFMGAN, AnomalyDiffusion and our method to generate anomaly image-mask pairs. These pairs are used to train BiSeNet V2 (Yu et al., 2021), UPerNet (Xiao et al., 2018) and LFD (Zhou et al., 2024a) respectively. Different from AnomalyDiffusion (Hu et al., 2024), which trains one segmentation model per product, we train a unified segmentation model for all the products. We also compare the segmentation results based on SeaS with the state-of-the-art unsupervised anomaly detection methods, such as RealNet (Zhang et al., 2024) and HVQ-Trans (Lu et al., 2023), in appendix A.5.

**Anomaly image generation quality.** In Tab. 1, we compare SeaS with some state-of-the-art AG and AIG methods on generation fidelity (IS) and diversity (IC-LPIPS). SeaS achieves 1.88 on IS and 0.34 on IC-LPIPS, which demonstrates that our method generates anomaly images with higher fidelity and diversity. We exhibit the generated anomaly images in Fig. 5, the anomaly images

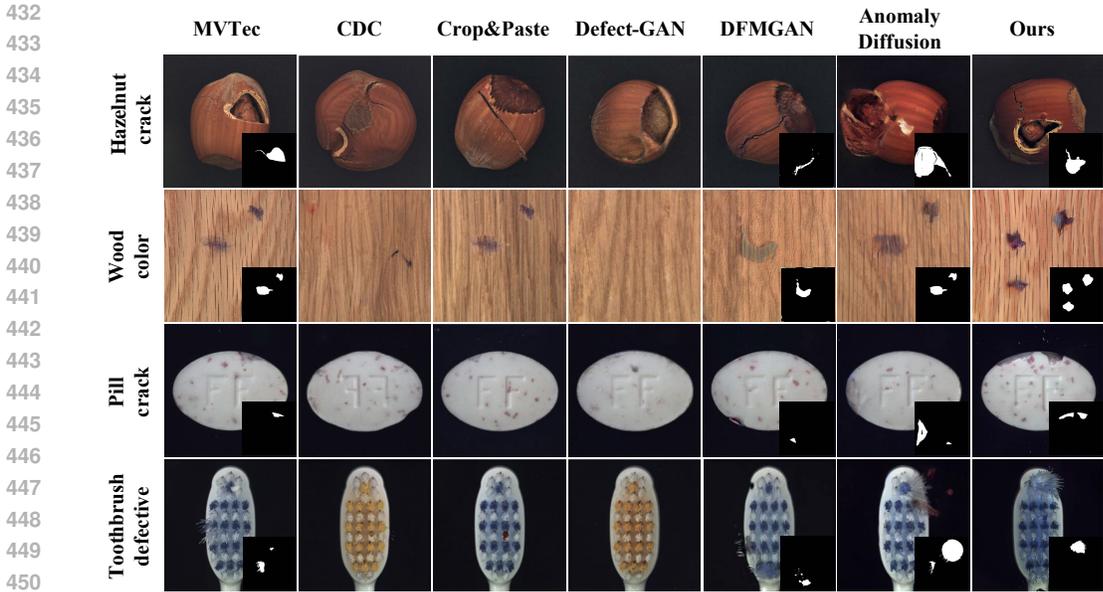


Figure 5: Visualization of the generation results on MVTec AD. The sub-image in the lower right corner is the generated mask, none means that the method cannot generate masks.

Table 2: Comparison on anomaly segmentation on MVTec AD.

Model	DRAEM				DFMGAN				AnomalyDiffusion				Ours			
	AUROC	AP	$F_1$ -max	IoU	AUROC	AP	$F_1$ -max	IoU	AUROC	AP	$F_1$ -max	IoU	AUROC	AP	$F_1$ -max	IoU
BiSeNet V2	81.37	38.90	42.62	39.39	94.57	60.42	60.54	45.83	96.27	64.50	62.27	42.89	97.21	69.21	66.37	55.28
UPerNet	83.21	42.78	45.97	42.03	92.33	57.01	56.91	46.64	96.87	69.92	66.95	50.80	97.87	74.42	70.70	61.24
LFD	76.41	40.99	43.91	35.61	94.91	67.06	65.09	45.49	96.30	69.77	66.99	45.77	98.09	77.15	72.52	56.47
Average	80.33	40.89	44.17	39.01	93.94	61.50	60.85	45.99	96.48	68.06	65.40	46.49	<b>97.72</b>	<b>73.59</b>	<b>69.86</b>	<b>57.66</b>

Table 3: Comparison on image-level anomaly detection on MVTec AD.

Model	DRAEM			DFMGAN			AnomalyDiffusion			Ours		
	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max
BiSeNet V2	89.87	93.51	89.97	90.90	94.43	90.33	90.08	94.84	91.84	96.00	98.14	95.43
UPerNet	89.45	93.92	89.66	90.74	94.43	90.37	96.62	98.61	96.21	98.29	99.20	97.34
LFD	87.94	93.41	88.65	91.08	95.40	90.58	95.15	97.78	94.66	95.88	97.89	95.15
Average	89.09	93.61	89.43	90.91	94.75	90.43	93.95	97.08	94.24	<b>96.72</b>	<b>98.41</b>	<b>95.97</b>

generated by our method 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., *wood\_color* and *pill\_crack*). The masks generated by our method are also precisely aligned with the anomaly regions (e.g., *toothbrush\_defective*). More qualitative and quantitative anomaly image generation results are in appendix A.6.

**Anomaly segmentation and detection.** We generate 1,000 image-mask pairs for each anomaly type, and use the image-mask pairs of all products along with all the training normal images to train the unified segmentation model, rather than training separate segmentation models for each product. We test the models on the rest images of the testing set of MVTec AD, which are not included in the training set for generation. The results are given in Tab. 2. All the methods are trained using the same number of images and the same training settings, detailed in appendix A.7. The segmentation results consistently demonstrate that our method outperforms others across all the segmentation models, with an 11.17% average improvement on IoU. We show the segmentation anomaly maps in Fig. 6. By 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* and *carpet\_cut*. In addition, we use the maximum value of the segmentation anomaly map as the image-level anomaly score for anomaly detection. We report the image-level metrics in Tab. 3, and our method achieves a 2.77% gain on image-AUROC. More qualitative comparison results on anomaly segmentation are in appendix A.9 and appendix A.10.

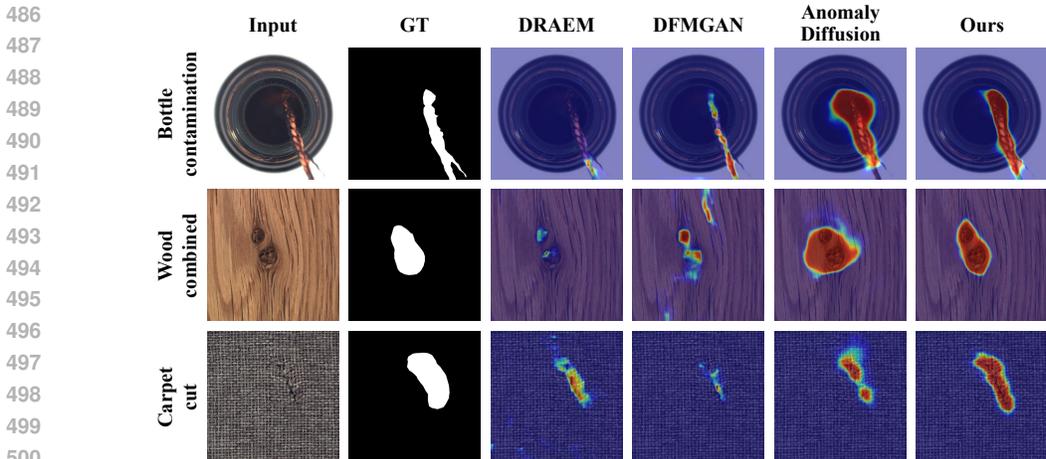


Figure 6: Qualitative anomaly segmentation results with BiSeNet V2 on MVTec AD.

### 4.3 ABLATION STUDY

**Anomaly image generation.** We train additional models to assess the effect of each component: **(a) the model with typical text prompt with fixed generic semantic words (short for with TP in Tab. 4); (b) the model without mixing the 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. 4); (e) our complete model.** We use these models to generate 1,000 anomaly image-mask pairs per anomaly type and train BiSeNet V2 for anomaly segmentation. In Tab. 4, the results show that omitting any component leads to a decrease in the fidelity and diversity of the generated images, as well as a decrease in the segmentation results. These validate the effectiveness of the components we proposed. More ablation studies on SeaS are shown in appendix A.8.

**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. 5) and coarse mask supervision (short for CMS in Tab. 5). **1) the model without any components, 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 the MRM in a progressive manner to refine coarse features; 4) our complete model.** We report the BiSeNet V2 results in Tab. 5, which demonstrates that each component in the RMP is indispensable for downstream anomaly segmentation. More ablation studies about RMP are in appendix A.8.

Table 4: 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)	1.88	0.34	97.21	69.21	66.37	55.28

Table 5: Ablation on the RMP branch.

Method	Metrics			
	MRM	PM	CMS	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
✓	✓	✓	✓	97.21
✓	✓	✓	✓	69.21
✓	✓	✓	✓	66.37
✓	✓	✓	✓	55.28

## 5 CONCLUSION

In this paper, we propose a novel few-shot industrial anomaly image generation method named SeaS. We explore an implicit characteristic that the anomalies exhibit randomness in shape and appearance, while the products maintain global consistency with minor variations in local details. We design a Separation and Sharing Fine-tuning strategy for industrial anomaly image generation, and a Refined Mask Prediction branch to obtain a fine-grained mask. Our method surpasses existing methods on both AIG and downstream anomaly segmentation tasks.

## REFERENCES

- 540  
541  
542 Omri Avrahami, Kfir Aberman, Ohad Fried, Daniel Cohen-Or, and Dani Lischinski. Break-a-scene:  
543 Extracting multiple concepts from a single image. In *SIGGRAPH Asia 2023 Conference Papers*,  
544 pp. 1–12, 2023.
- 545  
546 Shane Barratt and Rishi Sharma. A note on the inception score. *arXiv preprint arXiv:1801.01973*,  
547 2018.
- 548  
549 Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. Mvtec ad—a comprehen-  
550 sive real-world dataset for unsupervised anomaly detection. In *Proceedings of the IEEE/CVF*  
*Conference on Computer Vision and Pattern Recognition*, pp. 9592–9600, 2019.
- 551  
552 Paul Bergmann., Xin Jin., David Sattlegger., and Carsten Steger. The mvtec 3d-ad dataset for un-  
553 supervised 3d anomaly detection and localization. In *Proceedings of the 17th International Joint*  
*Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications*, pp.  
554 202–213, 2022.
- 555  
556 Mikołaj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd  
557 gans. In *International Conference on Learning Representations*, 2018.
- 558  
559 Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image  
560 editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
*Recognition*, pp. 18392–18402, 2023.
- 561  
562 Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite:  
563 Attention-based semantic guidance for text-to-image diffusion models. *ACM Transactions on*  
*Graphics*, 42(4):1–10, 2023.
- 564  
565 Hong Chen, Yipeng Zhang, Xin Wang, Xuguang Duan, Yuwei Zhou, and Wenwu Zhu. Disenbooth:  
566 Disentangled parameter-efficient tuning for subject-driven text-to-image generation. In *Interna-*  
*tional Conference on Learning Representations*, 2024.
- 567  
568 Yu-Min Chu, Chieh Liu, Ting-I Hsieh, Hwann-Tzong Chen, and Tyng-Luh Liu. Shape-guided dual-  
569 memory learning for 3d anomaly detection. In *Proceedings of the International Conference on*  
*Machine Learning*, pp. 6185–6194, 2023.
- 570  
571 Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks  
572 with cutout. *arXiv preprint arXiv:1708.04552*, 2017.
- 573  
574 Yuxuan Duan, Yan Hong, Li Niu, and Liqing Zhang. Few-shot defect image generation via defect-  
575 aware feature manipulation. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
576 volume 37, pp. 571–578, 2023.
- 577  
578 Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and  
579 Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using  
580 textual inversion. In *International Conference on Learning Representations*, 2022.
- 581  
582 Guan Gui, Bin-Bin Gao, Jun Liu, Chengjie Wang, and Yunsheng Wu. Few-shot anomaly-driven  
583 generation for anomaly classification and segmentation. In *European Conference on Computer*  
*Vision (ECCV 2024)*, pp. –, 2024.
- 584  
585 Ligong Han, Yinxiao Li, Han Zhang, Peyman Milanfar, Dimitris Metaxas, and Feng Yang. Svdiff:  
586 Compact parameter space for diffusion fine-tuning. In *Proceedings of the IEEE/CVF Interna-*  
*tional Conference on Computer Vision*, pp. 7323–7334, 2023.
- 587  
588 Haoyang He, Jiangning Zhang, Hongxu Chen, Xuhai Chen, Zhishan Li, Xu Chen, Yabiao Wang,  
589 Chengjie Wang, and Lei Xie. A diffusion-based framework for multi-class anomaly detection. In  
590 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 8472–8480, 2024.
- 591  
592 Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or.  
593 Prompt-to-prompt image editing with cross-attention control. In *International Conference on*  
*Learning Representations*, 2022.

- 594 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.  
595 Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in*  
596 *Neural Information Processing Systems*, 30, 2017.
- 597  
598 Teng Hu, Jiangning Zhang, Ran Yi, Yuzhen Du, Xu Chen, Liang Liu, Yabiao Wang, and Chengjie  
599 Wang. Anomalydiffusion: Few-shot anomaly image generation with diffusion model. In *Pro-*  
600 *ceedings of the AAAI Conference on Artificial Intelligence*, 2024.
- 601  
602 Jongheon Jeong, Yang Zou, Taewan Kim, Dongqing Zhang, Avinash Ravichandran, and Onkar  
603 Dabeer. Winclip: Zero-/few-shot anomaly classification and segmentation. In *Proceedings of the*  
604 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19606–19616, 2023.
- 605  
606 Chen Jin, Ryutaro Tanno, Amrutha Saseendran, Tom Diethe, and Philip Teare. An image is worth  
607 multiple words: Learning object level concepts using multi-concept prompt learning. In *Interna-*  
608 *tional Conference on Machine Learning*, 2024.
- 609  
610 Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept  
611 customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF Conference on Com-*  
612 *puter Vision and Pattern Recognition*, pp. 1931–1941, 2023.
- 613  
614 Chun-Liang Li, Kihyuk Sohn, Jinsung Yoon, and Tomas Pfister. Cutpaste: Self-supervised learning  
615 for anomaly detection and localization. In *Proceedings of the IEEE/CVF conference on Computer*  
616 *Vision and Pattern Recognition*, pp. 9664–9674, 2021.
- 617  
618 Daiqing Li, Huan Ling, Seung Wook Kim, Karsten Kreis, Sanja Fidler, and Antonio Torralba. Big-  
619 datasetgan: Synthesizing imagenet with pixel-wise annotations. In *Proceedings of the IEEE/CVF*  
620 *Conference on Computer Vision and Pattern Recognition*, pp. 21330–21340, 2022.
- 621  
622 Dongyun Lin, Yanpeng Cao, Wenbin Zhu, and Yiqun Li. Few-shot defect segmentation leveraging  
623 abundant defect-free training samples through normal background regularization and crop-and-  
624 paste operation. In *2021 IEEE International Conference on Multimedia and Expo*, pp. 1–6, 2021.
- 625  
626 Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense  
627 object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pp.  
628 2980–2988, 2017.
- 629  
630 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-*  
631 *ence on Learning Representations*, 2018.
- 632  
633 Ruiying Lu, YuJie Wu, Long Tian, Dongsheng Wang, Bo Chen, Xiyang Liu, and Ruimin Hu. Hier-  
634 archical vector quantized transformer for multi-class unsupervised anomaly detection. *Advances*  
635 *in Neural Information Processing Systems*, 36:8487–8500, 2023.
- 636  
637 Quang Nguyen, Truong Vu, Anh Tran, and Khoi Nguyen. Dataset diffusion: Diffusion-based syn-  
638 thetic data generation for pixel-level semantic segmentation. *Advances in Neural Information*  
639 *Processing Systems*, 36, 2024.
- 640  
641 Shuanlong Niu, Bin Li, Xinggang Wang, and Hui Lin. Defect image sample generation with gan  
642 for improving defect recognition. *IEEE Transactions on Automation Science and Engineering*,  
643 17(3):1611–1622, 2020.
- 644  
645 Utkarsh Ojha, Yijun Li, Jingwan Lu, Alexei A. Efros, Yong Jae Lee, Eli Shechtman, and Richard  
646 Zhang. Few-shot image generation via cross-domain correspondence. In *Proceedings of the*  
647 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10743–10752, 2021.
- 648  
649 Patrick Pérez, Michel Gangnet, and Andrew Blake. Poisson image editing. In *ACM SIGGRAPH*  
650 *2003*, pp. 313–318. 2003.
- 651  
652 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
653 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Con-*  
654 *ference on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.

- 648 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.  
649 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Pro-*  
650 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22500–  
651 22510, 2023.
- 652 Hannah M Schlüter, Jeremy Tan, Benjamin Hou, and Bernhard Kainz. Natural synthetic anomalies  
653 for self-supervised anomaly detection and localization. In *European Conference on Computer*  
654 *Vision*, pp. 474–489. Springer, 2022.
- 655 Weijia Wu, Yuzhong Zhao, Hao Chen, Yuchao Gu, Rui Zhao, Yefei He, Hong Zhou, Mike Zheng  
656 Shou, and Chunhua Shen. Datasetdm: Synthesizing data with perception annotations using diffu-  
657 sion models. *Advances in Neural Information Processing Systems*, 36:54683–54695, 2023a.
- 658 Weijia Wu, Yuzhong Zhao, Mike Zheng Shou, Hong Zhou, and Chunhua Shen. Diffumask: Syn-  
659 thesizing images with pixel-level annotations for semantic segmentation using diffusion models.  
660 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1206–1217,  
661 2023b.
- 662 Guangxuan Xiao, Tianwei Yin, William T Freeman, Frédo Durand, and Song Han. Fastcom-  
663 poser: Tuning-free multi-subject image generation with localized attention. *arXiv preprint*  
664 *arXiv:2305.10431*, 2023.
- 665 Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing  
666 for scene understanding. In *Proceedings of the European Conference on Computer Vision*, pp.  
667 418–434, 2018.
- 668 Jinheng Xie, Yuexiang Li, Yawen Huang, Haozhe Liu, Wentian Zhang, Yefeng Zheng, and  
669 Mike Zheng Shou. Boxdiff: Text-to-image synthesis with training-free box-constrained diffusion.  
670 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7452–7461,  
671 2023.
- 672 Xincheng Yao, Ruoqi Li, Zefeng Qian, Yan Luo, and Chongyang Zhang. Focus the discrepancy:  
673 Intra-and inter-correlation learning for image anomaly detection. In *Proceedings of the IEEE/CVF*  
674 *International Conference on Computer Vision*, pp. 6803–6813, 2023a.
- 675 Xincheng Yao, Ruoqi Li, Jing Zhang, Jun Sun, and Chongyang Zhang. Explicit boundary guided  
676 semi-push-pull contrastive learning for supervised anomaly detection. In *Proceedings of the*  
677 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 24490–24499, 2023b.
- 678 Changqian Yu, Changxin Gao, Jingbo Wang, Gang Yu, Chunhua Shen, and Nong Sang. Bisenet  
679 v2: Bilateral network with guided aggregation for real-time semantic segmentation. *International*  
680 *Journal of Computer Vision*, 129:3051–3068, 2021.
- 681 Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. Draem-a discriminatively trained reconstruc-  
682 tion embedding for surface anomaly detection. In *Proceedings of the IEEE/CVF International*  
683 *Conference on Computer Vision*, pp. 8330–8339, 2021.
- 684 Gongjie Zhang, Kaiwen Cui, Tzu-Yi Hung, and Shijian Lu. Defect-gan: High-fidelity defect syn-  
685 thesis for automated defect inspection. In *Proceedings of the IEEE/CVF Winter Conference on*  
686 *Applications of Computer Vision*, pp. 2524–2534, 2021a.
- 687 Hui Zhang, Zuxuan Wu, Zheng Wang, Zhineng Chen, and Yu-Gang Jiang. Prototypical residual  
688 networks for anomaly detection and localization. In *Proceedings of the IEEE/CVF Conference on*  
689 *Computer Vision and Pattern Recognition*, pp. 16281–16291, 2023a.
- 690 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image  
691 diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,  
692 pp. 3836–3847, 2023b.
- 693 Ximiao Zhang, Min Xu, and Xiuzhuang Zhou. Realnet: A feature selection network with realis-  
694 tic synthetic anomaly for anomaly detection. In *Proceedings of the IEEE/CVF Conference on*  
695 *Computer Vision and Pattern Recognition*, pp. 16699–16708, 2024.

702 Yuxuan Zhang, Huan Ling, Jun Gao, Kangxue Yin, Jean-Francois Lafleche, Adela Barriuso, Antonio  
703 Torralba, and Sanja Fidler. Datasetgan: Efficient labeled data factory with minimal human effort.  
704 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
705 10145–10155, 2021b.

706 Huan Zhou, Feng Xue, Yucong Li, Shi Gong, Yiqun Li, and Yu Zhou. Exploiting low-level rep-  
707 resentations for ultra-fast road segmentation. *IEEE Transactions on Intelligent Transportation*  
708 *Systems*, 2024a.

709 Qihang Zhou, Guansong Pang, Yu Tian, Shibo He, and Jiming Chen. Anomalyclip: Object-agnostic  
710 prompt learning for zero-shot anomaly detection. In *International Conference on Learning Rep-*  
711 *resentations*, 2024b.

712 Yang Zou, Jongheon Jeong, Latha Pemula, Dongqing Zhang, and Onkar Dabeer. Spot-the-difference  
713 self-supervised pre-training for anomaly detection and segmentation. In *European Conference on*  
714 *Computer Vision*, pp. 392–408. Springer, 2022.

715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A APPENDIX

### A.1 OVERVIEW

This supplementary material consists of:

- Analysis on decoupled anomaly alignment loss and multiple tokens (Sec. A.2).
- More implementation details (Sec. A.3).
- Anomaly image and mask generation results, downstream segmentation results on MVTeC 3D AD dataset (RGB images) (Sec. A.4).
- More quantitative comparison with unsupervised anomaly detection methods (Sec. A.5).
- More qualitative and quantitative results of anomaly image generation (Sec. A.6).
- More details of downstream segmentation model implementation and usage (Sec. A.7).
- More ablation studies (Sec. A.8), including ablation studies on the Unbalanced Abnormal Text Prompt design, the Separation and Sharing Fine-tuning loss, the minimum size requirement for training images, the training strategy of SeaS, the cross-attention maps for Decoupled Anomaly Alignment, the features for Coarse Feature Extraction, the features of VAE for Refined Mask Prediction, the normal image supervision for Refined Mask Prediction, the Mask Refinement Module, and the threshold for mask binarization.
- Qualitative comparison results of segmentation models trained on image-mask pairs generated by different anomaly generation methods (Sec. A.9).
- Qualitative comparison results of different segmentation models trained on image-mask pairs generated by SeaS (Sec. A.10).
- **Anomaly image and mask generation results, downstream segmentation results on VisA dataset (Sec. A.11).**
- **Explanation of discriminative features in U-Net decoder (Sec. A.12).**
- **Comparison with the Textual Inversion (Sec. A.13).**
- **More experiments on lighting conditions (Sec. A.14).**
- **More experiments on replacing generation strategies (Sec. A.15).**
- **More visualization on recombining the decoupled attributes for unseen anomalies (Sec. A.16).**

### A.2 ANALYSIS ON DECOUPLED ANOMALY ALIGNMENT LOSS AND MULTIPLE TOKENS

Here we give a more detailed analysis of the learning process of the DA loss. According to Eq. 3, intuitively, the DA loss may pull the anomaly tokens similar to each other. However, the U-Net in Stable Diffusion uses multi-head attention, which ensures different anomaly tokens cover different attributes of the anomalies. In Eq. 3, the cross-attention map is the multiply of the feature map of U-Net and the anomaly tokens. In the implementation of multi-head attention, both the learnable embedding of the anomaly token and the U-Net feature are decomposed into several groups along the channel dimension. E.g., the conditioning vector  $e_a \in \mathbb{R}^{1 \times C_1}$ , which is corresponding to anomaly token, is divided into  $\{e_{a,i} \in \mathbb{R}^{1 \times \frac{C_1}{q}} \mid i \in [1, q]\}$ , and the image feature  $v \in \mathbb{R}^{r \times r \times C_2}$  is divided into  $\{v_i \in \mathbb{R}^{1 \times \frac{C_2}{q}} \mid i \in [1, q]\}$ , where  $q$  is the number of heads in the multi-head attention. Then the corresponding groups are multiplied, and the outputs of all the heads are averaged. The attention map  $A$  of  $e_a$  is calculated by:

$$A = \frac{1}{q} \sum_{i=1}^q \text{softmax}\left(\frac{Q_i K_{a,i}^\top}{\sqrt{d}}\right), Q_i = \phi_q(v_i), K_{a,i} = \phi_k(e_{a,i}). \quad (8)$$

Therefore, in the defect region, the DA loss only ensures the average of each head tends to 1, but does not require the anomaly tokens to be the same with each other. In addition, each  $e_a$  is different from each other, and is combined by  $e_{a,i}$ . **The update direction of each  $e_{a,i}$  is related to  $v_i$  and covers some features of the defect, it encompasses the attributes of anomalies from various perspectives, thereby providing diversified information.**

### A.3 MORE IMPLEMENTATION DETAILS

**More training details.** For the Unbalanced Abnormal Text Prompt, we set the number  $N$  of multiple  $\langle df_n \rangle$  to 4 and the number  $N'$  of  $\langle ob \rangle$  to 1, these parameters are fixed across all product classes. For example, for the normal token  $\langle ob \rangle$ , given the lookup  $\mathcal{U} \in \mathbb{R}^{b \times 768}$ , where  $b$  is the number of text embeddings stored by the pre-trained text encoder, we use a placeholder string "ob1" as the input. Firstly, "ob1" is converted to a token ID  $s_{ob1} \in \mathbb{R}^{1 \times 1}$  in the tokenizer. Secondly,  $s_{ob1} \in \mathbb{R}^{1 \times 1}$  is converted to a one-hot vector  $\mathcal{S}_{ob1} \in \mathbb{R}^{1 \times (b+1)}$ . Thirdly, one learnable new embedding  $g \in \mathbb{R}^{1 \times 768}$  corresponding to  $s_{ob1}$  is inserted to the lookup  $\mathcal{U}$ , resulting in  $\mathcal{U}' \in \mathbb{R}^{(b+1) \times 768}$ . Here,  $g \in \mathbb{R}^{1 \times 768}$  is the learnable embedding of  $\langle ob \rangle$ . **These embeddings and U-Net are learnable during the fine-tuning process.**

**Training image generation model.** For each product, we perform  $800 \times G$  steps for fine-tuning, where  $G$  represents the number of anomaly categories of the product. The batch size of training image generation model is set to 4. During each step of our fine-tuning process, we sample 2 images from the abnormal training set  $X_{df}$ , and 2 images from the normal training set  $X_{ob}$ . We utilize the AdamW (Loshchilov & Hutter, 2018) optimizer with a learning rate of U-Net is  $4 \times 10^{-6}$ . The learning rate of the text embedding is  $4 \times 10^{-5}$ .

**Training Refined Mask Prediction branch.** We design a cascaded Refined Mask Prediction (RMP) branch, which is grafted onto the U-Net trained according to SeaS. For each product, we perform  $800 \times G$  steps for the RMP model, where  $G$  represents the number of anomaly types for the product. The batch size of training the RMP branch is set to 4. During each step of our fine-tuning process, we sample 2 images with their corresponding masks from the abnormal training set  $X_{df}$ , and 2 images from the normal training set  $X_{ob}$ . The masks used to suppress noise in normal images has each pixel value set to 0. The learning rate of the RMP model is  $5 \times 10^{-4}$ .

**Metrics.** For anomaly image generation, we report 2 metrics: the Inception Score (IS) and Intra-cluster pairwise LPIPS Distance (IC-LPIPS). The Inception Score (IS), proposed in (Barratt & Sharma, 2018), serves as an independent metric to evaluate the fidelity and diversity of generated images, by measuring the mutual information between input samples and their predicted classes. The IC-LPIPS (Ojha et al., 2021) is used to evaluate the diversity of generated images, which quantifies the perceptual similarity between image patches in the same cluster. For pixel-level anomaly segmentation and image-level anomaly detection, we report 3 metrics: Area Under Receiver Operator Characteristic curve (AUROC), Average Precision (AP), and  $F_1$ -score at the optimal threshold ( $F_1$ -max). **All of these metric are calculated using the *scikit-learn* library.** In addition, we calculate the Intersection over Union (IoU) to more accurately evaluate the anomaly segmentation result.

**Resource requirement and time consumption.** We conduct our training on NVIDIA Tesla A100 40G GPU. Specifically, we use a single A100 to train a generation model sequentially for each product category, with each training process occupying about 20G of GPU memory. Since each anomaly type requires isolate training, the training time depends on the total amount of anomaly types across all products. For example, the product *metal\_nut* contains 4 anomaly types, and each needs around 35 minutes. The generation model for *metal\_nut* spends 2 hours and 20 minutes on training in total. For the RMP branch, each anomaly type needs around 25 minutes. Hence, it takes 1 hour and 40 minutes to train *metal\_nut*. More details are given in Tab. 6, where  $K$  is the total number of anomaly types across all products. The comparison on time consumption is shown in Tab. 7. For the MVTecAD datasets with 73 anomaly types, our training takes 73 hours, which is shorter than the 249 hours required by AnomalyDiffusion and the 414 hours required by DFMGAN. In terms of inference time, SeaS costs 720 ms per image, which is shorter than the 3830 ms per image required by the Diffusion-based method AnomalyDiffusion. The inference time of the GAN-based method DFMGAN is 48ms per image.

Table 6: Computational resource and training time.

Stage	Time (minutes per product)	GPU(MB)	Overall Time
Generation Model	$35 \times K$	20242	42 hours and 35 minutes
RMP branch	$25 \times K$	23280	30 hours and 25 minutes

Table 7: Comparison on time consumption.

Model	Overall Training Time(hours)	Inference Time (ms)
DFMGAN (Duan et al., 2023)	414	48
AnomalyDiffusion (Hu et al., 2024)	249	3830
Ours	73	720

## A.4 ADDITIONAL DATASET RESULTS

We perform experimental evaluations on the RGB images of the MVTec 3D AD Dataset (Bergmann et al., 2022), which includes 10 product categories, each with up to 4 different anomalies. It encompasses several common challenges, such as variations in lighting conditions and product poses, which are crucial for validating the robustness of image generation methods. The experimental settings are the same as those in Sec. 4.1 and Sec. A.3.

Table 8: Comparison on IS and IC-LPIPS on MVTec 3D AD. Bold indicates the best performance.

Category	DFMGAN (Duan et al., 2023)		AnomalyDiffusion (Hu et al., 2024)		Ours	
	IS $\uparrow$	IC-L $\uparrow$	IS $\uparrow$	IC-L $\uparrow$	IS $\uparrow$	IC-L $\uparrow$
bagel	1.07	0.26	1.02	0.22	<b>1.28</b>	<b>0.29</b>
cable_gland	1.59	<b>0.25</b>	1.79	0.19	<b>2.21</b>	0.19
carrot	1.94	<b>0.29</b>	1.66	0.17	<b>2.07</b>	0.22
cookie	1.80	0.31	1.77	0.29	<b>2.07</b>	<b>0.38</b>
dowel	<b>1.96</b>	<b>0.37</b>	1.60	0.20	1.95	0.26
foam	1.50	0.17	1.77	0.30	<b>2.20</b>	<b>0.39</b>
peach	2.11	<b>0.34</b>	1.91	0.23	<b>2.40</b>	0.28
potato	<b>3.05</b>	<b>0.35</b>	1.92	0.17	1.98	0.22
rope	1.46	0.29	1.28	0.25	<b>1.53</b>	<b>0.41</b>
tire	1.53	0.25	1.35	0.20	<b>1.81</b>	<b>0.31</b>
Average	1.80	0.29	1.61	0.22	<b>1.95</b>	<b>0.30</b>

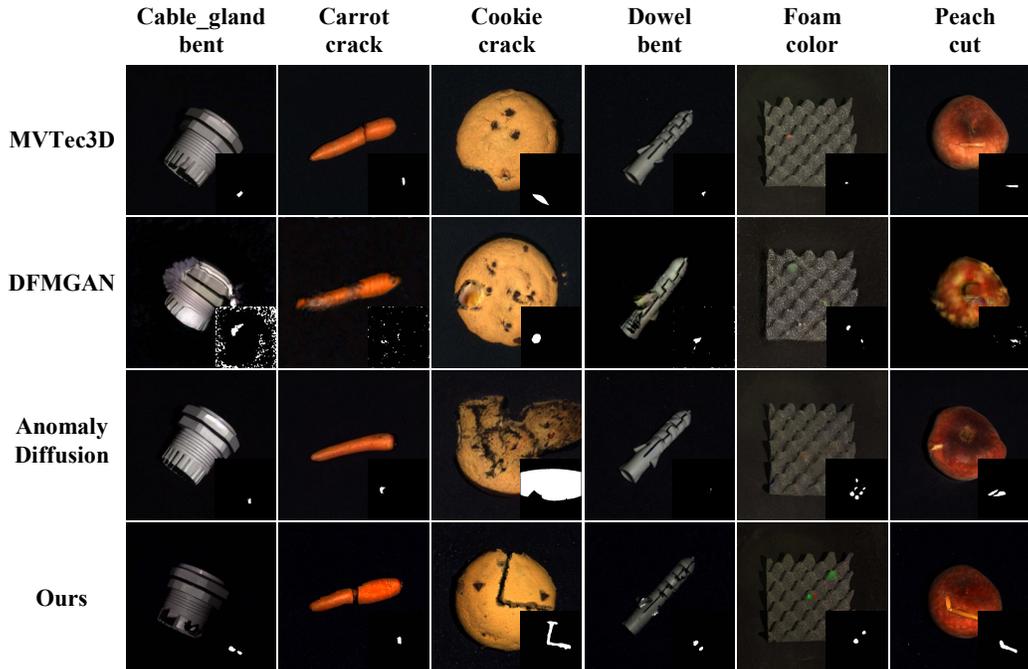


Figure 7: Visualization of the generation results on MVTec 3D AD. The sub-image in the lower right corner is the generated mask.

## Comparison Methods

In terms of compared approaches, since existing state-of-the-art approaches, e.g., DFMGAN(Duan et al., 2023), AnomalyDiffusion (Hu et al., 2024), conducted the experimental evaluations only on MVTec AD dataset (Bergmann et al., 2019), we evaluate them on MVTec 3D AD dataset (Bergmann. et al., 2022) using their official source codes.

### Anomaly image generation quality

As presented in Tab. 8, SeaS achieves scores of 1.95 on IS and 0.30 on IC-LPIPS, demonstrating our method’s ability to generate anomaly images with superior fidelity and diversity. The generated anomaly images are shown in Fig. 7. SeaS can generate images with diverse anomalies, avoiding overfitting to the training images (e.g., *cookie\_crack* and *foam\_color*), while ensuring the fidelity of the generated images (e.g., *cable\_gland\_bent*). Additionally, the masks generated by our method are accurately aligned with the anomalies (e.g., *peach\_cut*).

### Anomaly segmentation and detection

Tab. 9 shows the comparisons on downstream supervised segmentation trained by the generated images. It consistently demonstrates that our method outperforms others across all the segmentation models, with a 15.49% average improvement on IoU. The segmentation anomaly maps are shown in Fig. 8. There are fewer false positives (e.g., *potato\_combined*) and fewer false negatives (e.g., *bagel\_contamination*), when the BiSeNet V2 is trained on the image-mask pairs generated by our method. In addition, we report the image-level metrics in Tab. 10, and our method achieves a 6.74% gain on image-AUROC. Tab. 11 shows the comparisons of anomaly detection methods HVQ-Trans (Lu et al., 2023), Shape-guided (Chu et al., 2023), and FOD (Yao et al., 2023a) on anomaly segmentation tasks. Supervised segmentation models achieve better performance than most unsupervised AD methods on small-scale networks, with an IoU of 39.00% on LFD (0.936M). The pixel-level AUROC, which is sensitive to false negatives but less sensitive to false positives, of the Shape-guided method is higher. However, our observation indicates that the Shape-guided method has a high number of false positives. This significantly degrades the segmentation metrics, resulting in low pixel-level AP,  $F_1$ -max, and IoU scores. For industrial anomaly detection, an effective method should achieve a balance between false positives and false negatives.

Table 9: Comparison on anomaly segmentation on MVTec 3D AD.

Model	DFMGAN (Duan et al., 2023)				AnomalyDiffusion (Hu et al., 2024)				Ours			
	AUROC	AP	$F_1$ -max	IoU	AUROC	AP	$F_1$ -max	IoU	AUROC	AP	$F_1$ -max	IoU
BiSeNet V2 (Yu et al., 2021)	75.89	15.02	21.73	15.68	92.39	15.15	20.09	14.70	90.41	26.04	32.61	28.55
UPerNet (Xiao et al., 2018)	75.12	19.54	26.04	18.78	88.48	28.95	35.81	25.04	91.93	38.51	43.53	38.56
LFD (Zhou et al., 2024a)	72.15	9.54	14.29	14.81	92.68	24.29	32.74	19.90	91.61	40.25	43.47	39.00
Average	74.39	14.70	20.69	16.42	91.18	22.80	29.55	19.88	<b>91.32</b>	<b>34.93</b>	<b>39.87</b>	<b>35.37</b>

Table 10: Comparison on image-level anomaly detection on MVTec 3D AD.

Model	DFMGAN (Duan et al., 2023)			AnomalyDiffusion (Hu et al., 2024)			Ours		
	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max
BiSeNet V2 (Yu et al., 2021)	61.88	81.80	84.44	61.49	81.35	85.36	73.60	87.75	85.82
UPerNet (Xiao et al., 2018)	67.56	84.53	84.99	76.56	90.42	87.35	82.75	92.59	88.72
LFD (Zhou et al., 2024a)	62.23	82.17	85.38	77.06	89.44	87.20	78.96	91.22	87.28
Average	63.89	82.83	84.94	71.70	87.07	86.64	<b>78.44</b>	<b>90.52</b>	<b>87.27</b>

Table 11: Comparison with anomaly detection methods on MVTec 3D AD.

Model	Parameters	Image-level			Pixel-level			
		AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max	IoU
SeaS + BiSeNet V2 (Yu et al., 2021)	3.341M	73.60	87.75	85.82	90.41	26.04	32.61	28.55
SeaS + UPerNet (Xiao et al., 2018)	64.042M	<b>82.57</b>	<b>92.59</b>	<b>88.72</b>	91.93	38.51	<b>43.53</b>	38.56
SeaS + LFD (Zhou et al., 2024a)	<b>0.936M</b>	78.96	91.22	87.28	91.61	<b>40.25</b>	43.47	<b>39.00</b>
HVQ-Trans (Lu et al., 2023)	8.45M	68.15	84.38	85.20	96.40	24.59	17.23	20.51
Shape-guided (Chu et al., 2023)	4.13M	79.07	91.05	<b>88.72</b>	<b>98.45</b>	26.69	34.16	34.12
FOD (Yao et al., 2023a)	3.58M	71.66	86.83	86.57	97.03	14.70	20.99	23.31

972  
 973  
 974  
 975  
 976  
 977  
 978  
 979  
 980  
 981  
 982  
 983  
 984  
 985  
 986  
 987  
 988  
 989  
 990  
 991  
 992  
 993  
 994  
 995  
 996  
 997  
 998  
 999  
 1000  
 1001  
 1002  
 1003  
 1004  
 1005  
 1006  
 1007  
 1008  
 1009  
 1010  
 1011  
 1012  
 1013  
 1014  
 1015  
 1016  
 1017  
 1018  
 1019  
 1020  
 1021  
 1022  
 1023  
 1024  
 1025

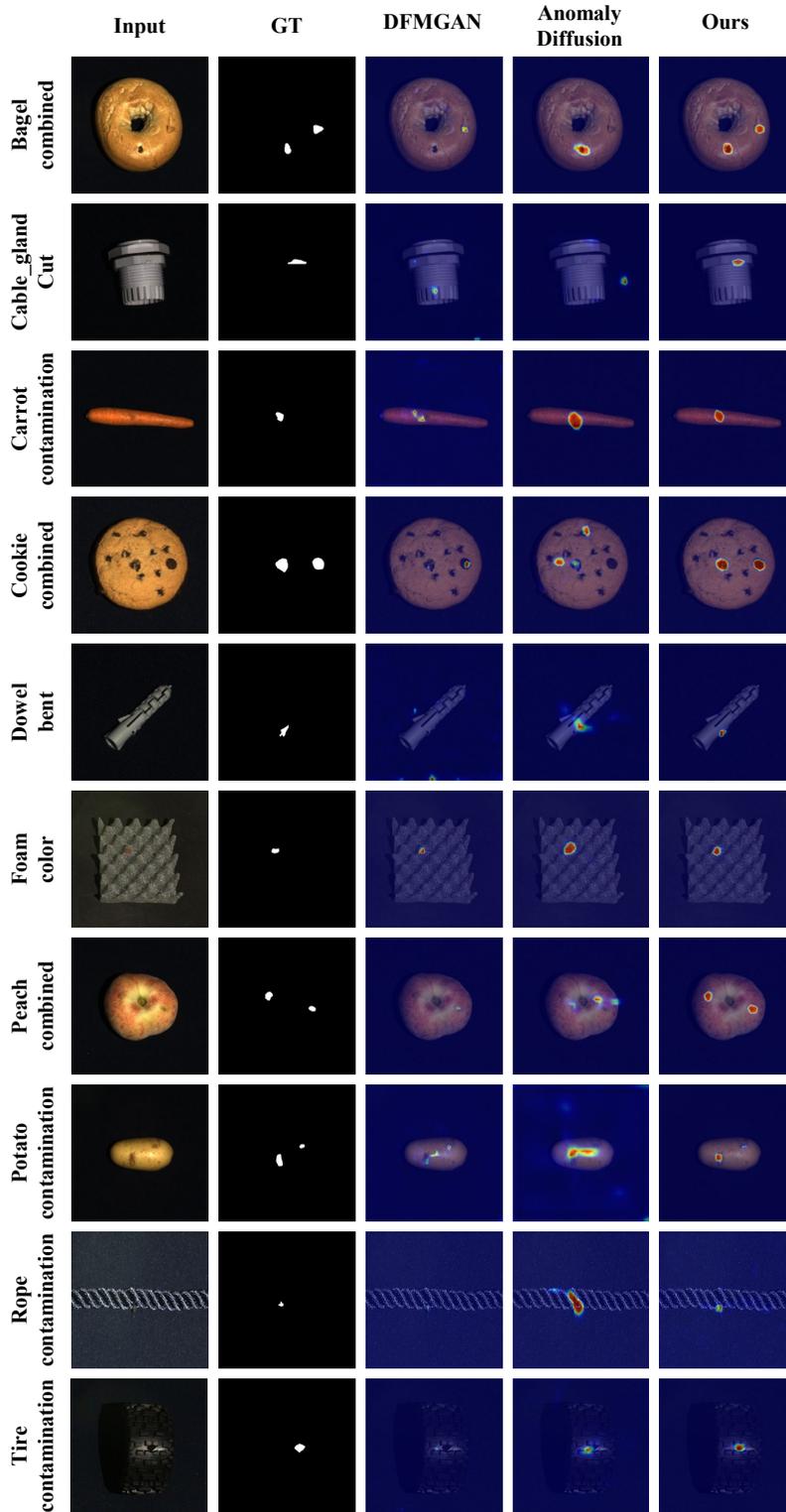


Figure 8: Qualitative anomaly segmentation results with BiSeNet V2 on MVTec 3D AD.

1026 **More qualitative anomaly image generation results**

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

We provide further qualitative results of every category on the MVTec 3D AD dataset, from Fig. 9 to Fig. 10. We report the anomaly image generation results of SeaS for varying types of anomalies. The first column represents the generated anomaly images, the second column represents the corresponding generated masks.



Figure 9: Qualitative results of our anomaly image generation results on MVTec 3D AD. In the first row, from left to right are the results for *bagel*, *cable\_gland*, and *carrot* categories. In the second row, from left to right are the results for *cookie*, *dowel*, and *foam* categories.

1080  
 1081  
 1082  
 1083  
 1084  
 1085  
 1086  
 1087  
 1088  
 1089  
 1090  
 1091  
 1092  
 1093  
 1094  
 1095  
 1096  
 1097  
 1098  
 1099  
 1100  
 1101  
 1102  
 1103  
 1104  
 1105  
 1106  
 1107  
 1108  
 1109  
 1110  
 1111  
 1112  
 1113  
 1114  
 1115  
 1116  
 1117  
 1118  
 1119  
 1120  
 1121  
 1122  
 1123  
 1124  
 1125  
 1126  
 1127  
 1128  
 1129  
 1130  
 1131  
 1132  
 1133

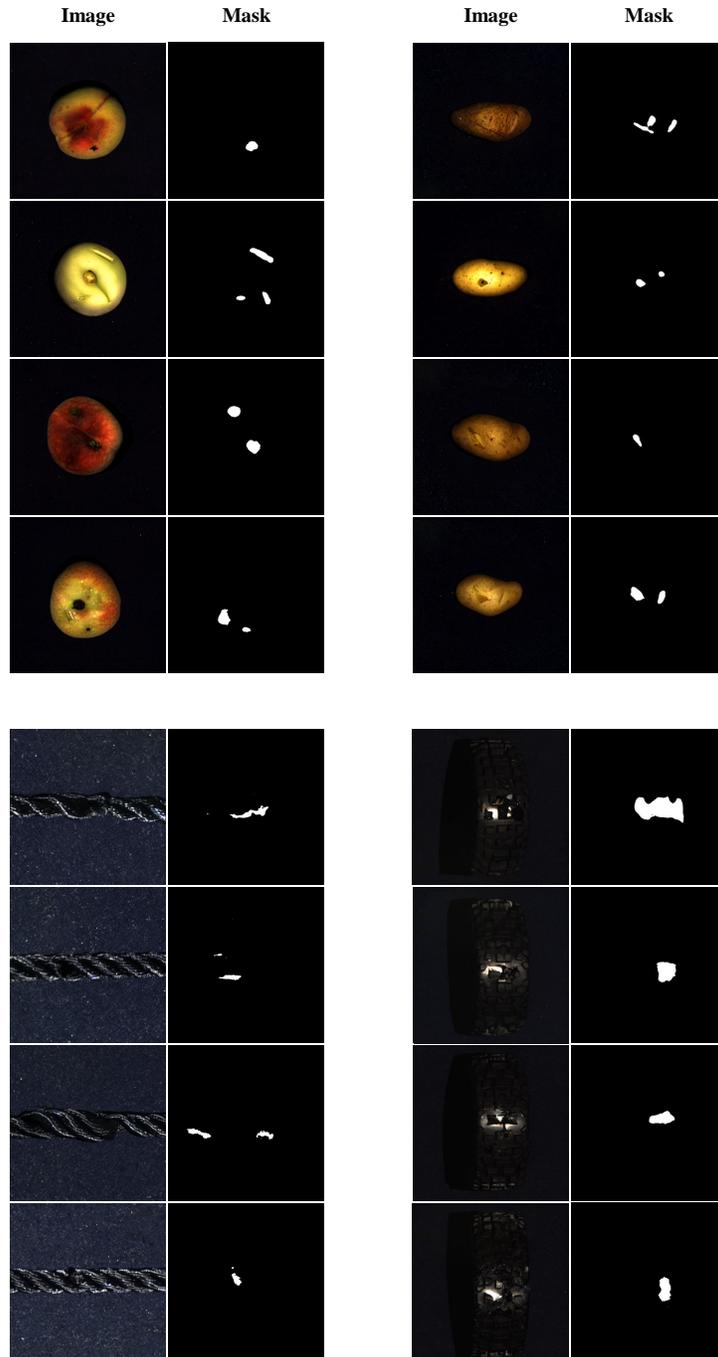


Figure 10: Qualitative results of our anomaly image generation results on MVTEC 3D AD. In the first row, from left to right are the results for *peach*, and *potato* categories. In the second row, from left to right are the results for *rope*, and *tire* categories.

**More qualitative segmentation results with different segmentation models**

In this section, we provide further qualitative results with the anomaly segmentation models on the MVTec 3D AD dataset. As shown in Fig. 11, we report the segmentation results of different segmentation models trained on image-mask pairs generated by SeaS.

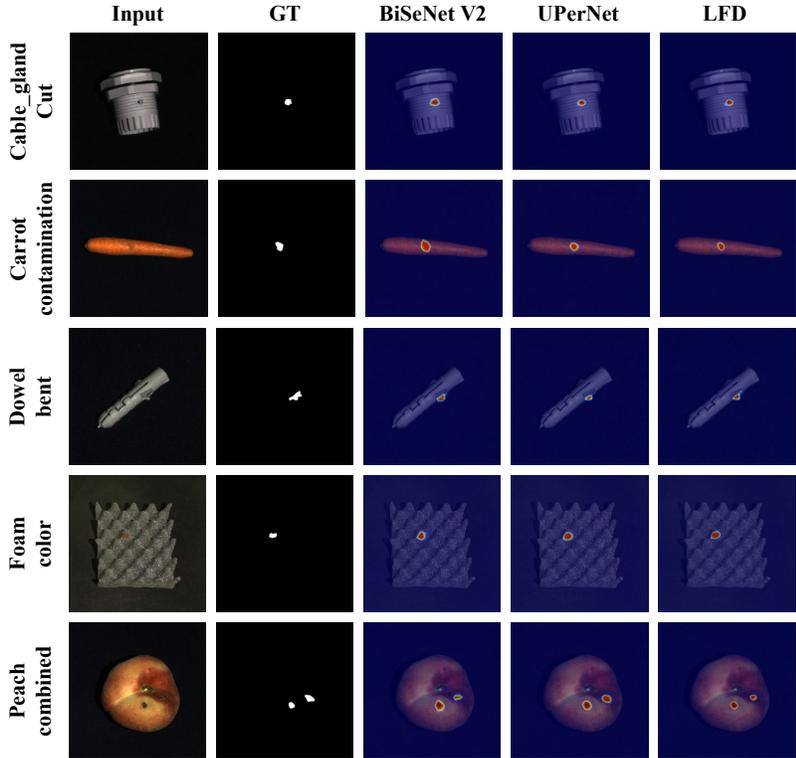


Figure 11: Qualitative comparison results with the anomaly segmentation models on MVTec 3D AD. In the figure, from top to bottom are the results for *cable\_gland*, *carrot*, *dowel*, *foam* and *peach* categories.

**A.5 MORE QUANTITATIVE COMPARISON WITH UNSUPERVISED ANOMALY DETECTION METHODS**

We compare our method with the state-of-the-art unsupervised anomaly detection methods, i.e., RealNet (Zhang et al., 2024), HVQ-Trans (Lu et al., 2023), DiAD (He et al., 2024), and PRN (Zhang et al., 2023a). The performance is different from the results reported in the paper, since as we mentioned in Sec 4.1, we use 2/3 anomaly images and all good images in the testing set of MVTec AD as the testing set in all the experiments, while the original results are achieved on the whole testing set. As shown in Tab. 12, Real-Net contains 591M parameters, around 177 times larger than BiSeNet V2, and 631 times larger than LFD, while the pixel-level AP and IoU measures of Real-Net are even worse than those of BiSeNet V2 and LFD. Although the pixel-level AUROC metric, which is more sensitive to false negatives than to false positives, is slightly higher for Real-Net, we observe that it generates a high number of false positives, substantially reducing pixel-level AP,  $F_1$ -max, and IoU scores. **The results of UperNet outperform the DiAD method on all measures, despite having only 1/24 of the parameters.** For effective industrial anomaly detection, a method must balance false positives and false negatives. The supervised anomaly segmentation models greatly outperform HVQ-Trans (8.45M parameters) and PRN in AP,  $F_1$ -max measure, and IoU, even though BiSeNet V2 and LFD are much smaller. These comparisons reveal that the small segmentation model achieves good performance using the generated images of SeaS, which is important for practical industrial applications.

Table 12: Comparison with anomaly detection methods on MVTec AD.

Model	Parameters	Image-level			Pixel-level			
		AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max	IoU
SeaS + BiSeNet V2 (Yu et al., 2021)	3.341M	96.00	98.14	95.43	97.21	69.21	66.37	55.28
SeaS + UPerNet (Xiao et al., 2018)	64.042M	<b>98.29</b>	<b>99.20</b>	97.34	97.87	74.42	70.70	<b>61.24</b>
SeaS + LFD (Zhou et al., 2024a)	<b>0.936M</b>	95.88	97.89	95.15	98.09	<b>77.15</b>	<b>72.52</b>	56.47
Real-Net (Zhang et al., 2024)	591M	98.19	98.99	<b>97.88</b>	<b>98.84</b>	68.09	66.46	53.99
HVQ-Trans (Lu et al., 2023)	8.45M	96.38	98.09	95.30	97.60	47.95	53.32	45.03
DiAD (He et al., 2024)	1525M	97.20	99.00	96.50	96.80	52.60	55.50	-
PRN (Zhang et al., 2023a)	-	91.60	96.60	92.40	96.90	66.20	64.70	-

## A.6 MORE QUALITATIVE AND QUANTITATIVE ANOMALY IMAGE GENERATION RESULTS

## More qualitative generation results

We provide further qualitative results of every category on the MVTec AD dataset, from Fig. 12 to Fig. 14. We report the anomaly image generation results of SeaS for varying types of anomalies. The first column represents the generated anomaly images, the second column represents the corresponding generated masks, and the third column represents the masks generated without using the Mask Refinement Module.

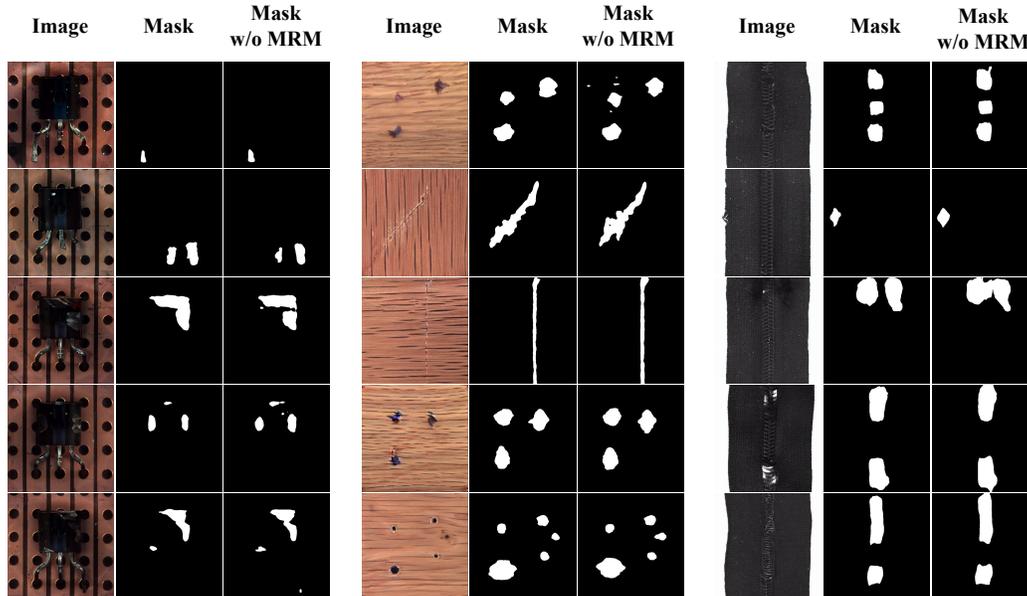


Figure 12: Qualitative results of our anomaly image generation results on MVTec AD. In the first row, from left to right are the results for *transistor*, *wood*, and *zipper* categories.

1242  
 1243  
 1244  
 1245  
 1246  
 1247  
 1248  
 1249  
 1250  
 1251  
 1252  
 1253  
 1254  
 1255  
 1256  
 1257  
 1258  
 1259  
 1260  
 1261  
 1262  
 1263  
 1264  
 1265  
 1266  
 1267  
 1268  
 1269  
 1270  
 1271  
 1272  
 1273  
 1274  
 1275  
 1276  
 1277  
 1278  
 1279  
 1280  
 1281  
 1282  
 1283  
 1284  
 1285  
 1286  
 1287  
 1288  
 1289  
 1290  
 1291  
 1292  
 1293  
 1294  
 1295

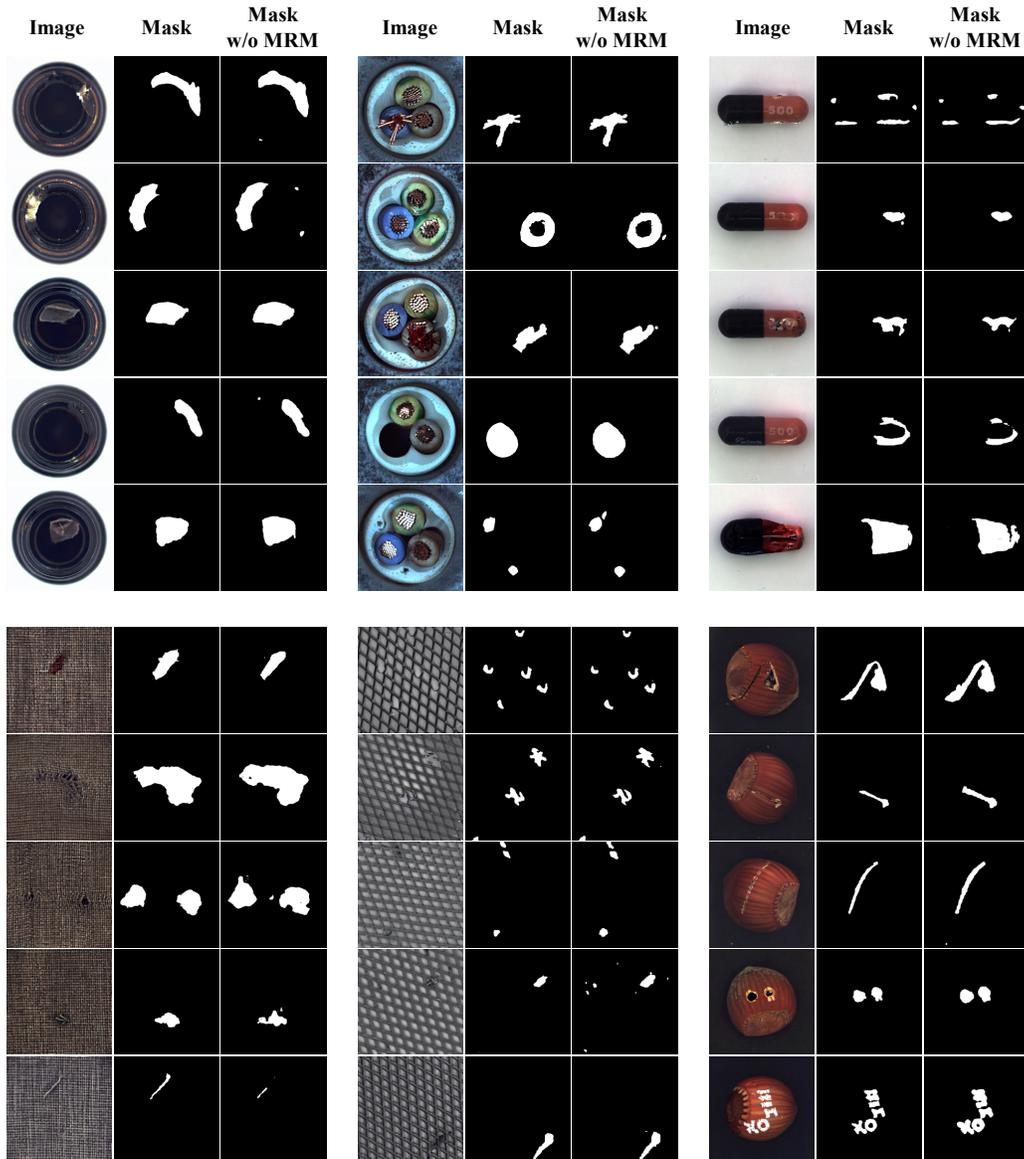


Figure 13: Qualitative results of our anomaly image generation results on MVTeC AD. In the first row, from left to right are the results for *bottle*, *cable*, and *capsule* categories. In the second row, from left to right are the results for *carpet*, *grid*, and *hazelnut* categories.

1296  
 1297  
 1298  
 1299  
 1300  
 1301  
 1302  
 1303  
 1304  
 1305  
 1306  
 1307  
 1308  
 1309  
 1310  
 1311  
 1312  
 1313  
 1314  
 1315  
 1316  
 1317  
 1318  
 1319  
 1320  
 1321  
 1322  
 1323  
 1324  
 1325  
 1326  
 1327  
 1328  
 1329  
 1330  
 1331  
 1332  
 1333  
 1334  
 1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
 1344  
 1345  
 1346  
 1347  
 1348  
 1349

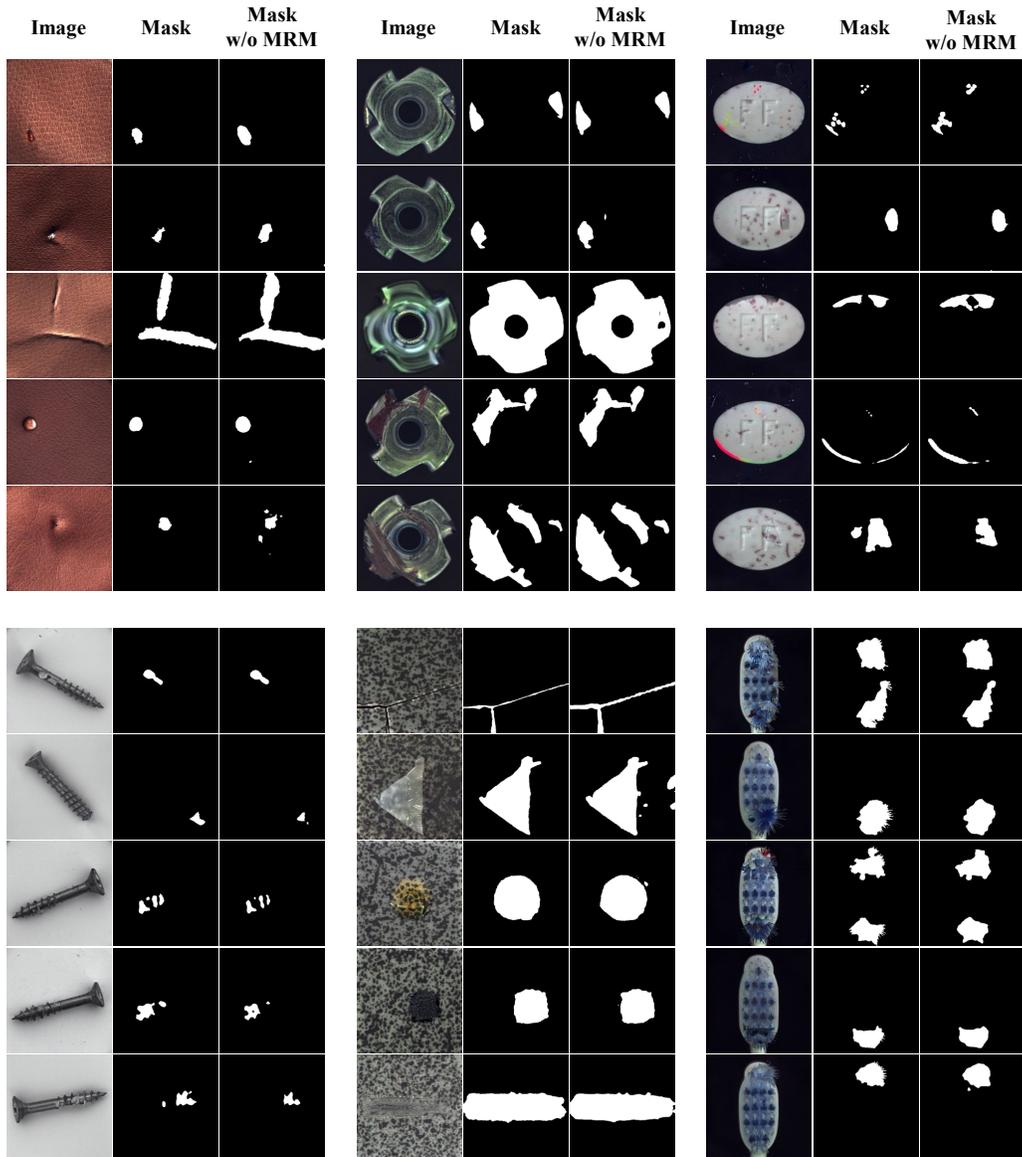


Figure 14: Qualitative results of our anomaly image generation results on MVTec AD. In the first row, from left to right are the results for *leather*, *metal\_nut*, and *pill* categories. In the second row, from left to right are the results for *screw*, *tile*, and *toothbrush* categories.

### More quantitative results

In this section, we report the detailed generation results of SeaS for each category on the MVTec AD datasets, compared with DRAEM (Zavrtanik et al., 2021), DFMGAN (Duan et al., 2023) and AnomalyDiffusion (Hu et al., 2024) which are presented from Tab. 13 to Tab. 18

Table 13: Comparison on anomaly segmentation on BiSeNet V2.

Category	DRAEM				DFMGAN				AnomalyDiffusion				Ours			
	AUROC	AP	$F_1$ -max	IoU	AUROC	AP	$F_1$ -max	IoU	AUROC	AP	$F_1$ -max	IoU	AUROC	AP	$F_1$ -max	IoU
bottle	77.42	45.25	48.21	45.77	89.34	64.67	62.78	44.71	99.00	88.02	80.53	68.25	<b>99.46</b>	<b>93.43</b>	<b>85.59</b>	<b>75.86</b>
cable	65.28	14.73	23.09	19.44	<b>93.87</b>	67.98	64.74	44.02	92.84	69.86	66.32	46.49	89.85	<b>72.07</b>	<b>71.58</b>	<b>53.24</b>
capsule	63.71	13.31	20.00	37.88	74.88	16.43	23.01	29.97	<b>92.71</b>	<b>38.11</b>	<b>40.67</b>	19.44	86.33	24.64	30.54	<b>39.70</b>
carpet	97.27	69.99	68.28	54.31	94.53	42.53	47.44	39.88	98.65	73.10	65.83	43.25	<b>99.61</b>	<b>82.30</b>	<b>72.94</b>	<b>55.52</b>
grid	93.86	25.36	35.46	35.57	96.86	24.40	37.40	29.93	80.59	8.08	16.79	14.26	<b>99.36</b>	<b>37.91</b>	<b>42.50</b>	<b>39.80</b>
hazelnut	77.48	41.52	48.36	54.58	<b>99.87</b>	<b>96.75</b>	<b>90.07</b>	<b>71.68</b>	97.71	63.34	59.87	43.12	97.82	78.55	73.09	68.47
leather	<b>99.76</b>	<b>67.02</b>	<b>64.96</b>	<b>52.37</b>	97.50	51.10	52.26	50.67	99.30	57.49	59.62	43.94	98.91	59.84	58.62	45.82
metal_nut	73.26	32.26	32.77	48.68	99.39	97.59	92.52	70.40	99.03	95.67	88.69	58.8	<b>99.69</b>	<b>98.29</b>	<b>93.23</b>	<b>74.40</b>
pill	60.02	9.33	17.17	11.67	97.09	83.98	79.26	36.39	<b>99.44</b>	<b>93.16</b>	<b>86.62</b>	41.18	98.31	76.97	68.00	<b>55.43</b>
screw	82.23	17.78	24.08	22.15	<b>97.94</b>	37.10	41.01	31.63	94.08	17.95	25.90	20.00	97.64	<b>40.20</b>	<b>45.35</b>	<b>38.43</b>
tile	98.09	82.9	76.43	63.48	99.65	97.08	91.16	<b>75.94</b>	97.79	85.58	78.28	60.46	<b>99.67</b>	<b>97.29</b>	<b>91.48</b>	75.75
toothbrush	92.65	36.73	45.92	23.90	97.70	<b>51.32</b>	54.05	23.38	<b>98.43</b>	49.64	<b>54.08</b>	26.53	97.15	46.09	49.02	<b>28.56</b>
transistor	62.48	14.83	20.57	21.85	84.31	45.34	46.07	30.00	<b>98.85</b>	<b>85.27</b>	<b>77.95</b>	49.83	96.75	69.52	66.11	<b>57.24</b>
wood	92.89	70.82	68.34	58.05	98.32	64.82	63.11	<b>58.99</b>	96.78	63.38	60.31	45.73	<b>98.38</b>	<b>80.81</b>	<b>74.03</b>	56.22
zipper	84.18	41.68	45.65	41.08	97.29	65.18	63.24	49.93	98.81	78.89	72.66	62.03	<b>99.23</b>	<b>80.27</b>	<b>73.41</b>	<b>64.80</b>
Average	81.37	38.90	42.62	39.39	94.57	60.42	60.54	45.83	96.27	64.5	62.27	42.89	<b>97.21</b>	<b>69.21</b>	<b>66.37</b>	<b>55.28</b>

Table 14: Comparison on image-level anomaly detection on BiSeNet V2.

Category	DRAEM			DFMGAN			AnomalyDiffusion			Ours		
	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max
bottle	95.93	98.21	93.18	96.74	98.75	95.35	98.14	99.34	97.67	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
cable	80.79	83.53	79.22	79.47	85.00	74.13	<b>95.37</b>	<b>96.71</b>	<b>92.91</b>	94.61	96.39	89.83
capsule	<b>91.88</b>	<b>97.62</b>	<b>92.41</b>	85.51	95.16	89.82	84.06	95.01	89.74	88.81	96.92	89.21
carpet	<b>98.21</b>	99.24	95.31	91.42	96.29	88.89	90.55	96.41	90.32	98.16	<b>99.31</b>	<b>97.56</b>
grid	96.43	98.50	95.00	<b>99.64</b>	<b>99.82</b>	97.56	81.19	89.92	83.95	99.17	99.63	<b>98.73</b>
hazelnut	97.92	98.65	94.62	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	93.39	95.74	90.91	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
leather	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	98.31	99.23	95.24	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	95.83	98.38	95.93
metal_nut	96.38	99.00	96.83	97.37	99.16	94.66	99.01	99.66	97.71	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
pill	74.68	89.88	89.00	84.86	95.27	91.00	90.38	97.43	91.35	<b>96.59</b>	<b>99.12</b>	<b>95.24</b>
screw	71.15	83.52	<b>83.15</b>	74.95	85.50	80.72	58.18	75.32	81.25	<b>77.24</b>	<b>89.55</b>	80.60
tile	99.68	99.82	98.28	99.47	99.74	99.12	98.78	99.44	97.39	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
toothbrush	82.50	90.00	85.11	78.33	87.73	83.72	78.33	89.26	79.17	<b>90.42</b>	<b>94.49</b>	<b>89.47</b>
transistor	73.87	68.65	60.87	79.52	75.77	69.57	94.40	94.68	94.34	<b>99.23</b>	<b>98.39</b>	<b>94.92</b>
wood	97.24	98.99	96.30	98.87	99.46	97.67	90.48	94.12	93.33	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
zipper	91.31	97.01	90.24	98.97	99.64	97.56	98.89	99.62	97.56	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
Average	89.87	93.51	89.97	90.90	94.43	90.33	90.08	94.84	91.84	<b>96.00</b>	<b>98.14</b>	<b>95.43</b>

Table 15: Comparison on anomaly segmentation on UPerNet.

Category	DRAEM				DFMGAN				AnomalyDiffusion				Ours			
	AUROC	AP	$F_1$ -max	IoU	AUROC	AP	$F_1$ -max	IoU	AUROC	AP	$F_1$ -max	IoU	AUROC	AP	$F_1$ -max	IoU
bottle	82.87	51.45	52.25	54.58	87.94	56.89	56.56	45.41	<b>99.54</b>	<b>93.01</b>	<b>85.94</b>	75.31	99.28	91.73	84.53	<b>78.73</b>
cable	56.64	13.59	22.42	16.29	87.52	64.30	65.61	41.02	91.00	68.12	67.49	51.84	<b>91.08</b>	<b>76.25</b>	<b>74.63</b>	<b>59.00</b>
capsule	60.95	11.12	18.34	27.18	67.92	12.31	20.32	30.47	<b>97.64</b>	<b>51.90</b>	<b>51.66</b>	37.00	92.09	39.60	43.89	<b>50.18</b>
carpet	98.83	79.57	72.53	61.50	95.85	36.05	34.52	48.10	99.45	<b>82.13</b>	72.55	53.17	<b>99.67</b>	82.01	<b>73.53</b>	<b>60.60</b>
grid	95.66	36.69	43.55	36.52	97.49	29.67	36.15	31.37	94.22	28.97	38.50	32.93	<b>99.18</b>	<b>44.94</b>	<b>48.28</b>	<b>44.21</b>
hazelnut	77.69	41.57	47.13	58.88	99.36	79.76	71.10	72.90	97.77	70.48	67.93	54.47	<b>99.54</b>	<b>81.84</b>	<b>75.48</b>	<b>73.30</b>
leather	<b>99.70</b>	65.49	63.01	<b>62.20</b>	80.97	17.60	26.21	30.17	99.48	63.46	60.54	48.70	99.42	<b>68.26</b>	<b>65.52</b>	57.01
metal_nut	65.37	23.26	27.39	42.56	98.44	95.64	91.48	64.92	98.62	95.11	88.62	61.31	<b>99.70</b>	<b>98.33</b>	<b>92.90</b>	<b>76.07</b>
pill	64.46	11.33	20.28	13.10	97.58	83.74	80.02	42.33	<b>99.33</b>	<b>95.04</b>	<b>88.77</b>	49.18	98.59	81.16	74.26	<b>62.62</b>
screw	90.88	23.64	31.49	24.71	97.49	<b>53.83</b>	<b>53.02</b>	42.05	93.89	36.60	42.68	34.08	<b>98.97</b>	52.02	51.65	<b>46.61</b>
tile	96.25	79.31	74.18	66.79	<b>99.79</b>	<b>97.29</b>	<b>91.11</b>	77.46	94.70	73.34	67.79	58.54	99.67	95.89	90.71	<b>77.89</b>
toothbrush	93.86	46.93	58.92	26.76	97.42	51.09	59.23	28.33	97.52	60.67	59.46	33.98	<b>98.50</b>	<b>63.62</b>	<b>63.07</b>	<b>42.09</b>
transistor	78.20	26.52	30.34	26.07	82.07	36.31	39.48	27.44	<b>94.26</b>	<b>73.68</b>	<b>69.50</b>	53.64	93.88	70.37	68.12	<b>56.98</b>
wood	95.03	77.07	74.07	64.67	97.90	69.02	62.21	63.10	96.09	70.10	64.38	51.44	<b>99.28</b>	<b>85.28</b>	<b>76.28</b>	<b>65.09</b>
zipper	91.74	54.11	53.69	48.57	97.28	71.60	66.64	54.54	<b>99.54</b>	<b>86.18</b>	<b>78.50</b>	66.47	99.17	85.01	77.57	<b>68.21</b>
Average	83.21	42.78	45.97	42.03	92.33	57.01	56.91	46.64	96.87	69.92	66.95	50.80	<b>97.87</b>	<b>74.42</b>	<b>70.70</b>	<b>61.24</b>

Table 16: Comparison on image-level anomaly detection on UPerNet.

Category	DRAEM			DFMGAN			AnomalyDiffusion			Ours		
	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max
bottle	98.49	99.41	97.62	94.19	97.86	93.18	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
cable	73.06	77.88	71.94	85.64	90.03	80.33	<b>95.58</b>	<b>97.06</b>	<b>92.56</b>	94.40	96.38	92.44
capsule	85.62	95.66	89.02	81.04	94.26	87.01	<b>96.00</b>	<b>98.77</b>	<b>95.48</b>	94.43	98.44	92.21
carpet	97.64	99.00	95.08	96.72	98.58	93.75	98.68	99.53	98.36	<b>99.94</b>	<b>99.97</b>	<b>99.20</b>
grid	97.62	98.99	97.44	98.33	99.13	96.30	96.67	98.73	97.44	<b>99.76</b>	<b>99.88</b>	<b>98.73</b>
hazelnut	97.14	97.74	92.63	99.84	99.87	97.96	99.17	99.43	97.87	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
leather	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	79.91	90.70	81.75	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
metal_nut	86.79	96.02	87.39	98.30	99.38	97.71	98.65	99.62	98.41	<b>99.72</b>	<b>99.91</b>	<b>99.21</b>
pill	67.07	87.37	88.07	88.54	96.56	92.39	91.23	97.78	90.91	<b>98.28</b>	<b>99.58</b>	<b>97.92</b>
screw	80.04	90.95	81.03	89.01	94.54	88.24	85.06	93.87	85.33	<b>93.47</b>	<b>97.07</b>	<b>90.45</b>
tile	99.20	99.51	98.25	99.68	99.81	99.13	99.68	99.81	99.13	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
toothbrush	86.67	94.19	88.89	75.00	86.99	80.00	90.00	95.13	90.00	<b>95.00</b>	<b>97.65</b>	<b>94.74</b>
transistor	79.29	74.68	69.09	83.04	73.59	74.19	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	99.52	99.16	96.43
wood	98.25	99.29	96.47	93.36	95.60	95.45	98.62	99.49	97.62	<b>99.87</b>	<b>99.94</b>	<b>98.82</b>
zipper	94.86	98.13	92.02	98.48	99.51	98.14	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
Average	89.45	93.92	89.66	90.74	94.43	90.37	96.62	98.61	96.21	<b>98.29</b>	<b>99.20</b>	<b>97.34</b>

Table 17: Comparison on anomaly segmentation on LFD.

Category	DRAEM				DFMGAN				AnomalyDiffusion				Ours			
	AUROC	AP	$F_1$ -max	IoU	AUROC	AP	$F_1$ -max	IoU	AUROC	AP	$F_1$ -max	IoU	AUROC	AP	$F_1$ -max	IoU
bottle	78.76	50.40	51.05	39.74	90.41	61.51	58.49	40.19	98.71	89.64	81.55	67.10	<b>99.28</b>	<b>92.65</b>	<b>84.86</b>	<b>73.82</b>
cable	67.64	19.41	25.36	19.19	96.49	79.40	<b>75.25</b>	53.47	<b>97.89</b>	<b>79.85</b>	72.75	53.69	94.53	75.41	72.70	<b>55.98</b>
capsule	88.48	34.60	39.62	31.48	91.82	<b>56.11</b>	<b>58.56</b>	32.50	<b>95.80</b>	38.17	48.92	32.04	91.80	49.76	53.69	<b>41.14</b>
carpet	83.51	37.69	41.86	41.22	89.10	48.04	49.89	39.46	94.83	53.15	51.79	42.21	<b>99.10</b>	<b>82.74</b>	<b>74.51</b>	<b>57.56</b>
grid	92.13	45.75	48.84	27.66	89.18	34.89	41.21	19.21	85.19	24.32	34.76	18.22	<b>98.78</b>	<b>62.24</b>	<b>58.44</b>	<b>41.69</b>
hazelnut	59.97	28.87	37.82	30.38	<b>99.36</b>	<b>95.16</b>	<b>89.80</b>	<b>76.43</b>	98.54	77.39	70.42	45.97	98.97	88.00	81.77	73.39
leather	97.38	66.36	63.43	53.13	97.82	51.86	52.25	48.09	98.99	65.73	62.85	42.65	<b>99.11</b>	<b>76.49</b>	<b>69.30</b>	<b>56.51</b>
metal_nut	63.34	35.42	35.24	47.51	98.16	95.16	90.99	63.02	<b>99.38</b>	<b>97.34</b>	<b>91.63</b>	64.59	99.23	96.66	91.42	<b>75.15</b>
pill	36.18	8.93	12.79	13.62	95.80	75.90	70.31	31.73	<b>98.96</b>	<b>92.51</b>	<b>85.35</b>	50.04	98.11	79.63	72.54	<b>56.73</b>
screw	91.03	27.05	32.95	19.03	93.96	38.00	41.69	30.88	92.68	44.64	49.17	34.08	<b>98.27</b>	<b>52.40</b>	<b>52.32</b>	<b>41.02</b>
tile	91.77	80.02	77.19	56.27	97.37	88.79	82.05	66.30	92.98	79.59	73.52	55.08	<b>99.38</b>	<b>96.24</b>	<b>89.90</b>	<b>75.50</b>
toothbrush	55.94	24.87	35.54	12.75	95.17	55.21	53.95	28.83	<b>98.31</b>	<b>68.60</b>	<b>66.14</b>	<b>29.67</b>	96.97	54.84	53.19	27.91
transistor	55.81	19.36	24.38	32.58	97.68	<b>89.68</b>	<b>84.18</b>	46.98	98.20	83.97	75.84	44.22	<b>98.80</b>	84.32	77.02	<b>55.57</b>
wood	90.04	73.42	71.25	59.55	97.47	77.72	70.91	58.77	95.68	67.54	63.06	42.78	<b>98.60</b>	<b>88.57</b>	<b>81.46</b>	<b>62.94</b>
zipper	94.15	62.65	61.39	50.12	93.80	58.43	56.82	46.44	98.42	84.05	77.08	64.14	<b>99.15</b>	<b>86.67</b>	<b>79.09</b>	<b>69.37</b>
Average	76.41	40.99	43.91	35.61	94.91	67.06	65.09	45.49	96.30	69.77	66.99	45.77	<b>98.01</b>	<b>77.77</b>	<b>72.81</b>	<b>57.62</b>

Table 18: Comparison on image-level anomaly detection on LFD.

Category	DRAEM			DFMGAN			AnomalyDiffusion			Ours		
	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max
bottle	96.40	98.56	94.12	96.98	98.76	95.35	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
cable	70.91	78.06	71.52	90.98	94.21	88.14	<b>99.52</b>	<b>99.55</b>	<b>97.71</b>	92.05	94.95	88.70
capsule	85.80	95.77	87.80	86.32	95.99	88.46	83.25	94.62	89.44	<b>93.80</b>	<b>98.19</b>	<b>93.42</b>
carpet	81.28	92.34	84.67	88.02	95.33	87.60	86.00	93.42	87.22	<b>97.98</b>	<b>99.22</b>	<b>96.67</b>
grid	<b>97.14</b>	98.63	92.86	85.48	92.61	85.71	93.69	97.08	91.14	96.79	<b>98.76</b>	<b>96.10</b>
hazelnut	85.73	89.00	81.72	99.90	99.91	98.97	98.28	98.60	95.83	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
leather	98.81	99.28	98.44	95.93	98.15	93.65	99.90	99.95	99.20	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
metal_nut	94.67	98.39	93.44	96.16	98.57	96.18	<b>99.01</b>	<b>99.65</b>	<b>98.46</b>	98.58	99.54	97.64
pill	66.39	89.29	88.07	82.85	94.40	92.00	94.15	98.42	94.47	<b>98.16</b>	<b>99.50</b>	<b>96.84</b>
screw	73.92	85.38	82.72	82.60	92.15	82.22	81.54	91.32	82.05	<b>87.83</b>	<b>94.39</b>	<b>85.54</b>
tile	95.80	97.49	93.33	98.94	99.43	96.55	98.25	99.13	95.65	<b>99.36</b>	<b>99.69</b>	<b>99.12</b>
toothbrush	83.33	90.29	84.44	77.08	87.68	80.95	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	87.92	94.08	87.80
transistor	90.06	89.06	81.48	88.04	85.06	77.78	97.38	96.57	92.86	<b>98.10</b>	<b>96.90</b>	<b>94.55</b>
wood	99.50	99.78	97.62	99.87	99.94	98.82	97.24	98.70	96.47	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
zipper	99.39	99.77	97.56	97.07	98.78	96.25	99.01	99.71	99.39	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
Average	87.94	93.41	88.65	91.08	95.40	90.58	95.15	97.78	94.66	<b>96.70</b>	<b>98.35</b>	<b>95.76</b>

## 1458 A.7 MORE DETAILS OF THE SEGMENTATION MODELS

1459  
1460 As mentioned in the experiment part, we choose three segmentation models (BiSeNet V2 (Yu et al.,  
1461 2021), UPerNet (Xiao et al., 2018), LFD (Zhou et al., 2024a)) to verify the validity of the generated  
1462 image-mask pairs on the downstream anomaly segmentation as well as detection tasks. **For BiSeNet  
1463 V2 and UPerNet, we generally follow the implementation provided by MMsegmentation. For  
1464 LFD, we also use the official implementation.**

1465 Specifically, for BiSeNet V2, we choose a backbone structure of a detail branch of three stages with  
1466 64, 64 and 128 channels and a semantic branch of four stages with 16, 32, 64 and 128 channels  
1467 respectively, with a decode head and four auxiliary heads (corresponding to the number of stages in  
1468 the semantic branch). As for UPerNet, we choose ResNet-50 as the backbone, with a decode head  
1469 and an auxiliary head.

1470 **In training segmentation models for downstream tasks, we adopt a training strategy of training  
1471 a unified segmentation model for all classes of products, rather than training separate segmen-  
1472 tation models for each class.** Experimental results are shown in Tab. 19, which indicate that the  
1473 performance of the unified segmentation model surpasses that of multiple individual segmentation  
1474 models.

1475 Table 19: Ablation on the training strategy of segmentation models.

Models	Multiple Models				Unified Model			
	AUROC	AP	$F_1$ -max	IoU	AUROC	AP	$F_1$ -max	IoU
BiSeNet V2	96.00	67.68	65.87	54.11	97.21	69.21	66.37	55.28
UPerNet	96.77	73.88	70.49	60.37	97.87	74.42	70.70	61.24
LFD	93.02	72.97	71.56	55.88	98.09	77.15	72.52	56.47
Average	95.26	71.51	69.31	56.79	<b>97.72</b>	<b>73.59</b>	<b>69.86</b>	<b>57.66</b>

## 1483 A.8 MORE ABLATION STUDIES

### 1485 Ablation on the Unbalanced Abnormal Text Prompt design

1486 In the design of the prompt for industrial anomaly image generation, we conduct experiments to  
1487 validate the effectiveness of our Unbalanced Abnormal (UA) Text Prompt for each anomaly type of  
1488 each product. We set the number of learnable  $\langle df_n \rangle$  to  $N$ , and the number of learnable  $\langle ob_j \rangle$  to  
1489  $N'$ . As shown in Tab. 20, by utilizing the UA Text Prompt, i.e.,

$$1490 \mathcal{P} = a \langle ob \rangle \text{ with } \langle df_1 \rangle, \langle df_2 \rangle, \langle df_3 \rangle, \langle df_4 \rangle$$

1491 we are able to provide high-fidelity and diverse images for downstream anomaly segmentation tasks,  
1492 resulting in the best performance in segmentation metrics.

### 1493 Ablation on the Separation and Sharing Fine-tuning loss

1494 In the design of the DA loss and NA loss for the Separation and Sharing Fine-tuning, we conduct  
1495 two sets of experiments: (a) We remove the second term in the DA loss (short for w/o ST in Tab.  
1496 21); (b) We replace the second term in DA loss with another term in the NA loss (short for with AT  
1497 in Tab. 21), which aligns the background area with the token  $\langle ob \rangle$  according to the mask:

$$1499 \mathcal{L}_{ob} = \sum_{l=1}^L (\|A_{ob}^l - (1 - M^l)\|^2) + \|\epsilon_{ob} - \epsilon_{\theta}(\hat{z}_{ob}, t_{ob}, \mathbf{e}_{ob})\|_2^2 \quad (9)$$

1500 where  $A_{ob}^l \in \mathbb{R}^{r \times r \times 1}$  is the cross-attention map corresponding to the normal token  $\langle ob \rangle$ . As  
1501 shown in Tab. 21, the experimental results demonstrate that, our adopted loss design achieves the  
1502 best performance in downstream segmentation tasks.  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511

Table 20: Ablation on the Unbalanced Abnormal Text Prompt design.

Prompt	AUROC	AP	$F_1$ -max	IoU
$N' = 1, N = 1$	96.48	63.69	62.50	52.02
$N' = 1, N = 4$ (Ours)	<b>97.21</b>	<b>69.21</b>	<b>66.37</b>	<b>55.28</b>
$N' = 4, N = 4$	96.55	66.28	63.95	54.07

Table 21: Ablation on the Separation and Sharing Fine-tuning loss.

Loss	AUROC	AP	$F_1$ -max	IoU
w/o ST	96.44	67.73	65.23	54.99
with AT	96.42	63.99	62.43	53.36
<b>Ours</b>	<b>97.21</b>	<b>69.21</b>	<b>66.37</b>	<b>55.28</b>

**Ablation on the minimum size requirement for training images**

In the few-shot setting, for a fair comparison, we follow the common setting in DFMGAN (Duan et al., 2023) and AnomalyDiffusion (Hu et al., 2024), i.e., using one-third abnormal image-mask pairs for each anomaly type in training. In this setting, the minimum number of abnormal training images is 2. Once we adopt a 3-shot setting, we need to reorganize the test set. To ensure that the test set is not reorganized for fair comparison, we take 1-shot and 2-shot settings for all anomaly types during training, i.e.,  $H = 1$  and  $H = 2$ , where  $H$  is the image number. The results are shown in Tab. 22 and Fig. 15. Observably, the models trained by 1-shot and 2-shot settings still generate anomaly images with decent diversity and authenticity.

Table 22: Ablation on the minimum size requirement for training images.

Size	IS	IC-L
$H = 1$	1.790	0.311
$H = 2$	1.794	0.314
$H = \frac{1}{3} \times H_0$	<b>1.876</b>	<b>0.339</b>

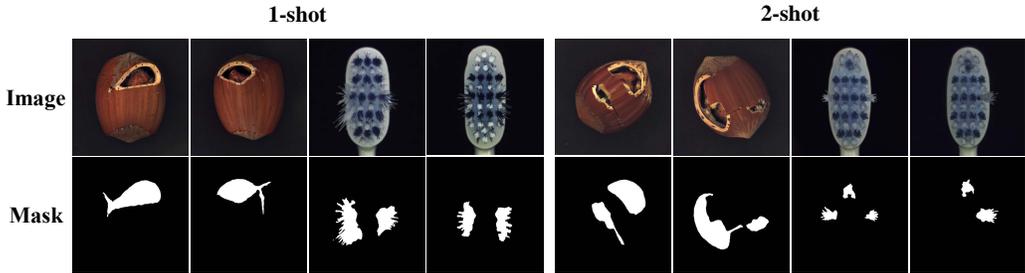


Figure 15: Visualization of the ablation study on the minimum size requirement for training images. In the figure, the first row is for generated images, the second row is for generated masks.

**Ablation on the training strategy of SeaS**

During each step of the fine-tuning process, we sample the same number of images from the abnormal training set  $X_{dr}$  and the normal training set  $X_{ob}$ . To investigate the efficacy of this strategy, we conduct three distinct sets of experiments: (a) prioritizing training with abnormal images followed by normal images (short for Abnormal-Normal in Tab. 23); (b) prioritizing training with abnormal images followed by anomaly images (short for Normal-Abnormal in Tab. 23); (c) training with a mix of both normal and abnormal images in each batch (short for Abnormal&Normal in Tab. 23). As shown in Tab. 23, SeaS yields superior performance in anomaly image generation, characterized by both high fidelity and diversity in the generated images.

**Ablation on the cross-attention maps for Decoupled Anomaly Alignment**

In Decoupled Anomaly Alignment (DA) loss, we leverage cross-attention maps from various layers of the U-Net encoder. Specifically, we investigate the impact of integrating different cross-attention maps, denoted as  $A^1 \in \mathbb{R}^{64 \times 64}$ ,  $A^2 \in \mathbb{R}^{32 \times 32}$ ,  $A^3 \in \mathbb{R}^{16 \times 16}$  and  $A^4 \in \mathbb{R}^{8 \times 8}$ . These correspond to the cross-attention maps of the “down-1”, “down-2”, “down-3”, and “down-4” layers of

Table 23: Ablation on training strategy of SeaS.

Strategy	IS	IC-L
Abnormal-Normal	1.53	0.28
Normal-Abnormal	1.70	0.32
Abnormal&Normal ( <b>Ours</b> )	<b>1.88</b>	<b>0.34</b>

the encoder in U-Net respectively. As shown in Tab. 24, the experimental results demonstrate that, employing a combination of  $\{A^2, A^3\}$  for DA loss, achieves the best performance in downstream segmentation tasks.

### Ablation on the features for Coarse Feature Extraction

In the coarse feature extraction process, we extract coarse but highly-discriminative features for anomalies from U-Net decoder. Specifically, we investigate the impact of integrating different features, denoted as  $F_1 \in \mathbb{R}^{16 \times 16 \times 1280}$ ,  $F_2 \in \mathbb{R}^{32 \times 32 \times 1280}$ ,  $F_3 \in \mathbb{R}^{64 \times 64 \times 640}$  and  $F_4 \in \mathbb{R}^{64 \times 64 \times 320}$ . These correspond to the output feature “up-1”, “up-2”, “up-3”, and “up-4” layers of the encoder in U-Net respectively. As shown in Tab. 25, the experimental results demonstrate that, employing a combination of  $\{F_2, F_3\}$  for coarse feature extraction, achieves the best performance in downstream segmentation task.

Table 24: Ablation on the cross-attention maps for Decoupled Anomaly Alignment.

$A^l$	AUROC	AP	$F_1$ -max	IoU
$l = 1, 2, 3$	96.42	68.92	66.24	54.52
$l = 2, 3, 4$	95.71	64.51	62.33	52.46
$l = 2, 3$ ( <b>Ours</b> )	<b>97.21</b>	<b>69.21</b>	<b>66.37</b>	<b>55.28</b>

Table 25: Ablation on the features for Coarse Feature Extraction.

$F_y$	AUROC	AP	$F_1$ -max	IoU
$y = 1, 2, 3$	94.35	63.58	60.54	52.36
$y = 2, 3, 4$	96.93	67.42	64.26	<b>55.31</b>
$y = 2, 3$ ( <b>Ours</b> )	<b>97.21</b>	<b>69.21</b>	<b>66.37</b>	55.28

### Ablation on the features of VAE for Refined Mask Prediction

In the Refined Mask Prediction, we combine the high-resolution features of VAE decoder features with discriminative features from U-Net, to generate accurately aligned anomaly image-mask pairs. In addition, we can also use the VAE encoder features as high-resolution features. As shown in Tab. 26, the experimental results show that, using VAE decoder features achieves better performance in downstream segmentation tasks.

Table 26: Ablation on the features of VAE for Refined Mask Prediction.

$F^{res}$	AUROC	AP	$F_1$ -max	IoU
VAE encoder	96.14	66.26	63.48	54.99
<b>VAE decoder</b>	<b>97.21</b>	<b>69.21</b>	<b>66.37</b>	<b>55.28</b>

### Ablation on the normal image supervision for Refined Mask Prediction

In the Refined Mask Prediction branch, we predict masks for normal images as the supervision for the mask prediction. We conduct two sets of experiments: (a) We remove the second and the fourth term in the loss for RMP, i.e., the normal image supervision (short for NIA in Tab. 27); (b) We use the complete form in RMP branch loss, i.e., we use the normal image for supervision, as in Eq. equation 10:

$$\mathcal{L}_{\mathcal{M}} = \mathcal{F}(\hat{M}_{df}, \mathbf{M}_{df}) + \mathcal{F}(\hat{M}_{ob}, \mathbf{M}_{ob}) + \mathcal{F}(\hat{M}'_{df}, \mathbf{M}'_{df}) + \mathcal{F}(\hat{M}'_{ob}, \mathbf{M}'_{ob}) \quad (10)$$

As shown in Tab. 27, the experimental results show that, using normal images for supervision achieves better performance in downstream segmentation tasks. We also provide further qualitative results of the effect of normal image supervision (short for NIA in Fig. 16) on MVTEC AD.

Table 27: Ablation on the normal image supervision for Refined Mask Prediction.

$F^{res}$	AUROC	AP	$F_1$ -max	IoU
w/o NIA	96.20	66.03	64.09	53.97
<b>with NIA (Ours)</b>	<b>97.21</b>	<b>69.21</b>	<b>66.37</b>	<b>55.28</b>

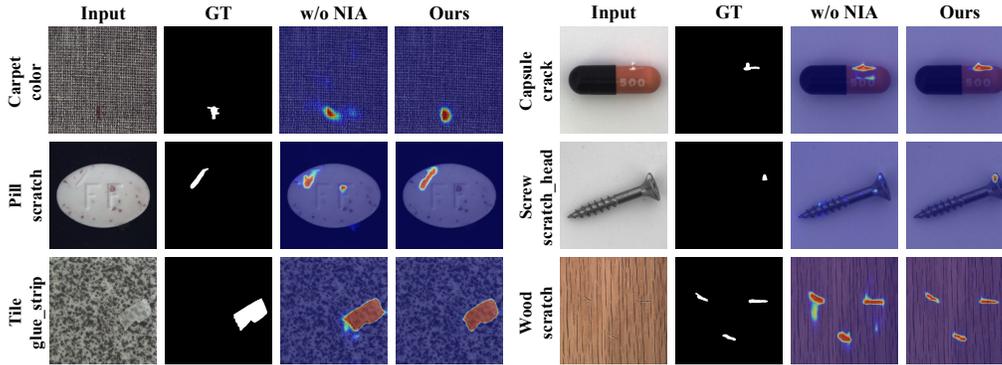


Figure 16: Qualitative results of the effect of normal image supervision on MVTec AD.

### Ablation on the Mask Refinement Module

In the Refined Mask Prediction branch, the Mask Refinement Module (MRM) is utilized to generate refined masks. We devise different structures for MRM, as shown in Fig. 17, including Case a): those without conv blocks, Case b): with one conv blocks, and Case c): with chained conv blocks. As shown in Fig. 18, we find that using the conv blocks in Case b), which consists of two  $1 \times 1$  convolutions and one  $3 \times 3$  convolution, helps the model learn the features of the defect area more accurately, rather than focusing on the background area for using one convolution alone in Case a). Based on this observation, we further designed a chained conv blocks structure in Case c), and the acquired features better reflect the defect area. This one-level-by-one level of residual learning helps the model achieve better residual correction results for the defect area features. As shown in Tab. 28 in the Appendix, Case c) improves the performance by + 0.28% on AUROC, + 2.29% on AP and + 2.29% on  $F_1$ -max, + 0.32% on IoU compared with Case b). We substantiate the superiority of the MRM structures that we design, through the results of downstream segmentation experiments.

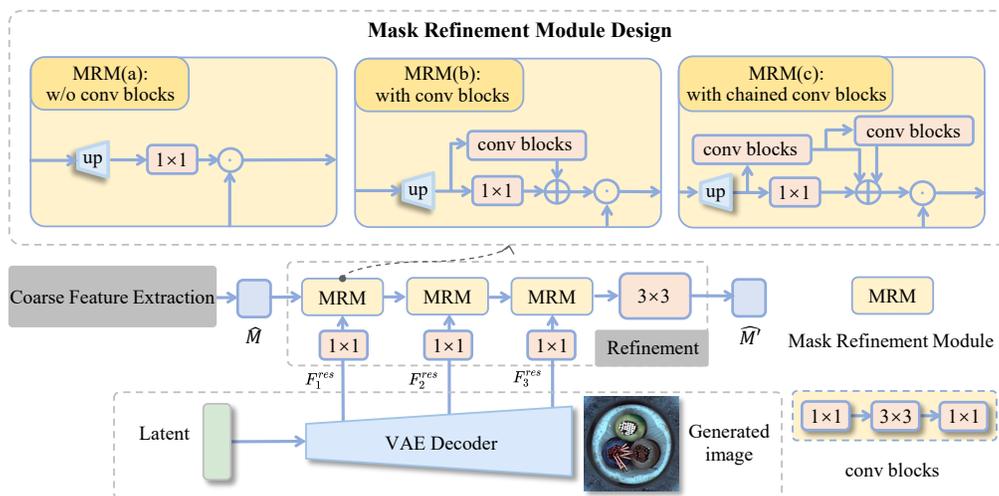


Figure 17: Different structure designs for the mask refinement module in the mask prediction branch.

1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727

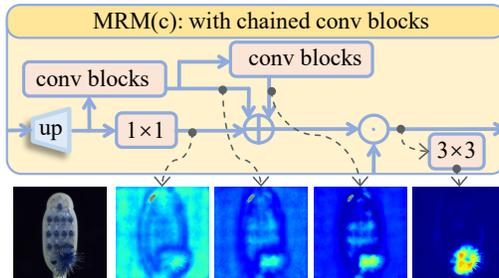


Figure 18: Visualization of the MRM module intermediate results. The top is for MRM structure diagram, and the bottom is sequentially for input image, feature maps of MRM intermediate process and the predicted mask.

Table 28: Ablation on the Mask Refinement Module.

Model	AUROC	AP	$F_1$ -max	IoU
with MRM (a)	96.75	68.18	64.96	<b>55.51</b>
with MRM (b)	96.93	66.92	64.08	54.96
<b>with MRM (c)</b>	<b>97.21</b>	<b>69.21</b>	<b>66.37</b>	55.28

#### Ablation on the threshold for mask binarization

In the Refined Mask Prediction branch, we take the threshold  $\tau$  for the second channel of refined anomaly masks  $\hat{M}'_{df}$  to segment the final anomaly mask. We train segmentation models using anomaly masks with  $\tau$  settings ranging from 0.1 to 0.5. As shown in Tab. 29, results indicate that setting  $\tau = 0.2$  yields the best model performance.

Table 29: Ablation on the threshold for mask binarization.

threshold	AUROC	AP	$F_1$ -max	IoU
$\tau = 0.1$	<b>97.56</b>	65.33	63.38	52.40
$\tau = 0.2$ (Ours)	97.21	<b>69.21</b>	<b>66.37</b>	<b>55.28</b>
$\tau = 0.3$	97.20	66.92	64.35	54.68
$\tau = 0.4$	95.31	63.55	61.97	53.03
$\tau = 0.5$	94.11	60.85	59.92	50.87

#### A.9 MORE QUALITATIVE COMPARISON RESULTS OF SEGMENTATION MODELS TRAINED ON IMAGE-MASK PAIRS GENERATED BY DIFFERENT ANOMALY GENERATION METHODS

We provide further qualitative results with different anomaly generation methods on the MVTEC AD dataset. We report the generation results of SeaS for varying types of anomalies in each category. Results are from Fig. 19 to Fig. 22.

1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781

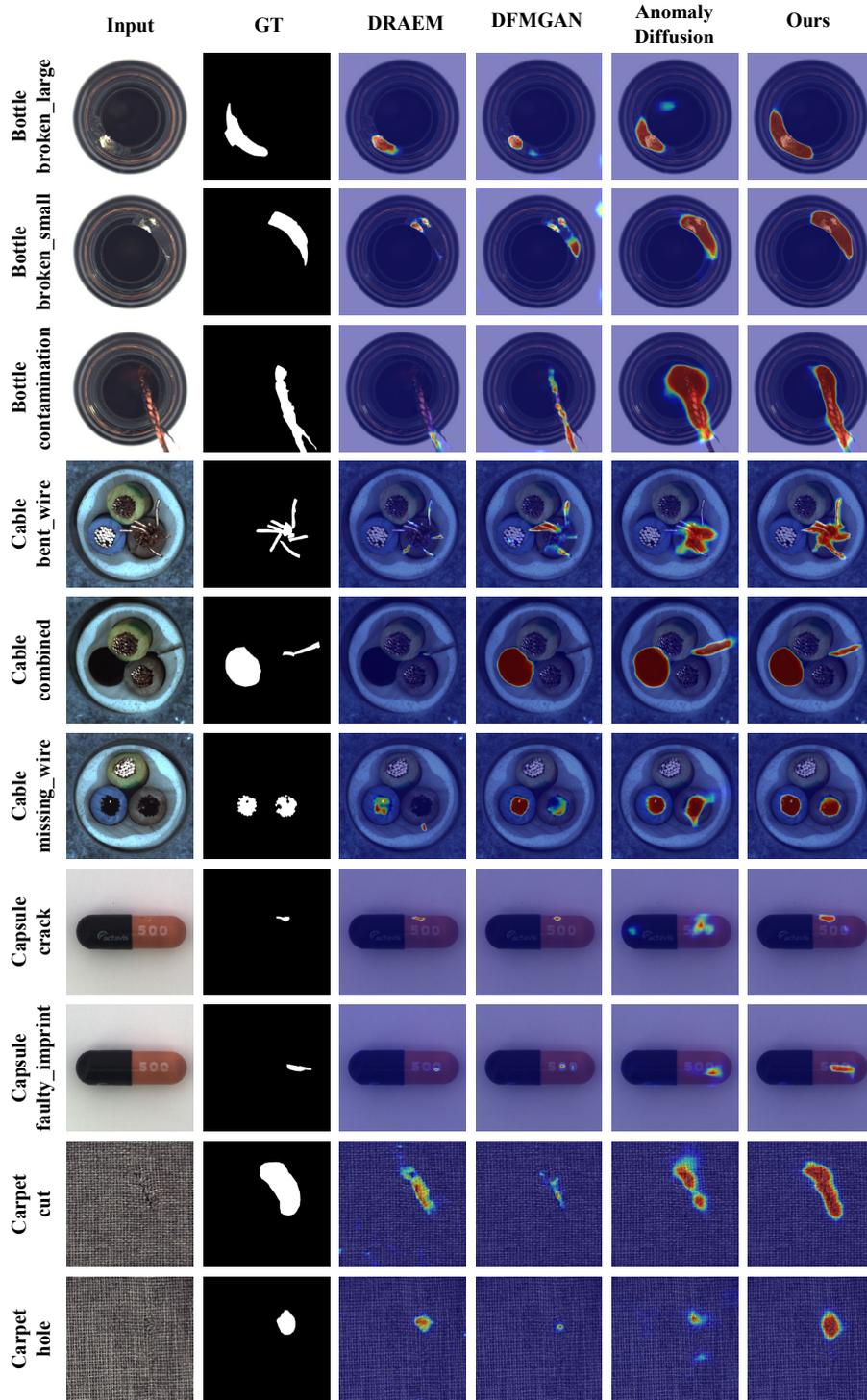


Figure 19: Comparison results with the anomaly segmentation models on MVTEC AD. In the figure, from top to bottom are the results for *bottle*, *cable*, *capsule* and *carpet* categories.

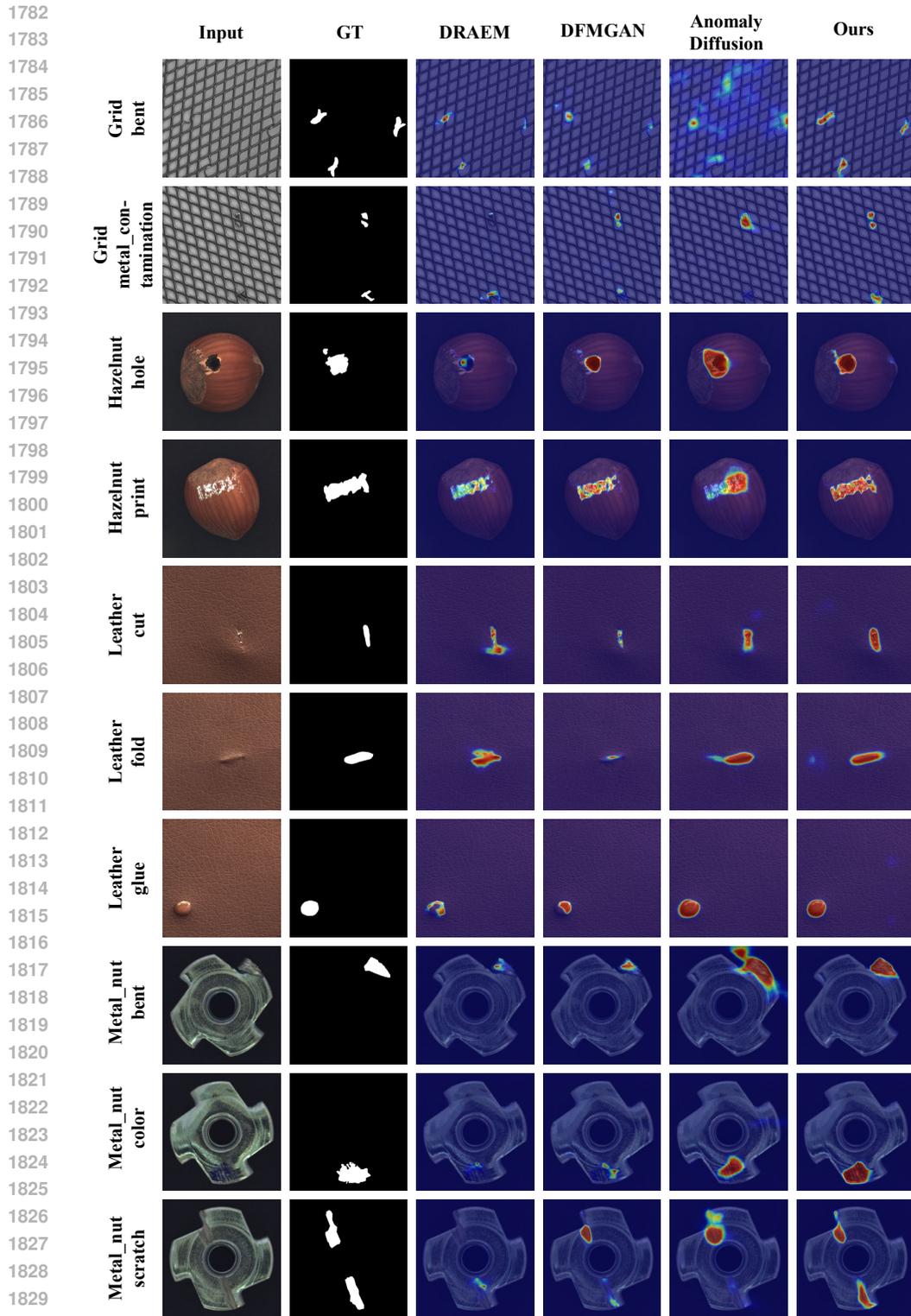
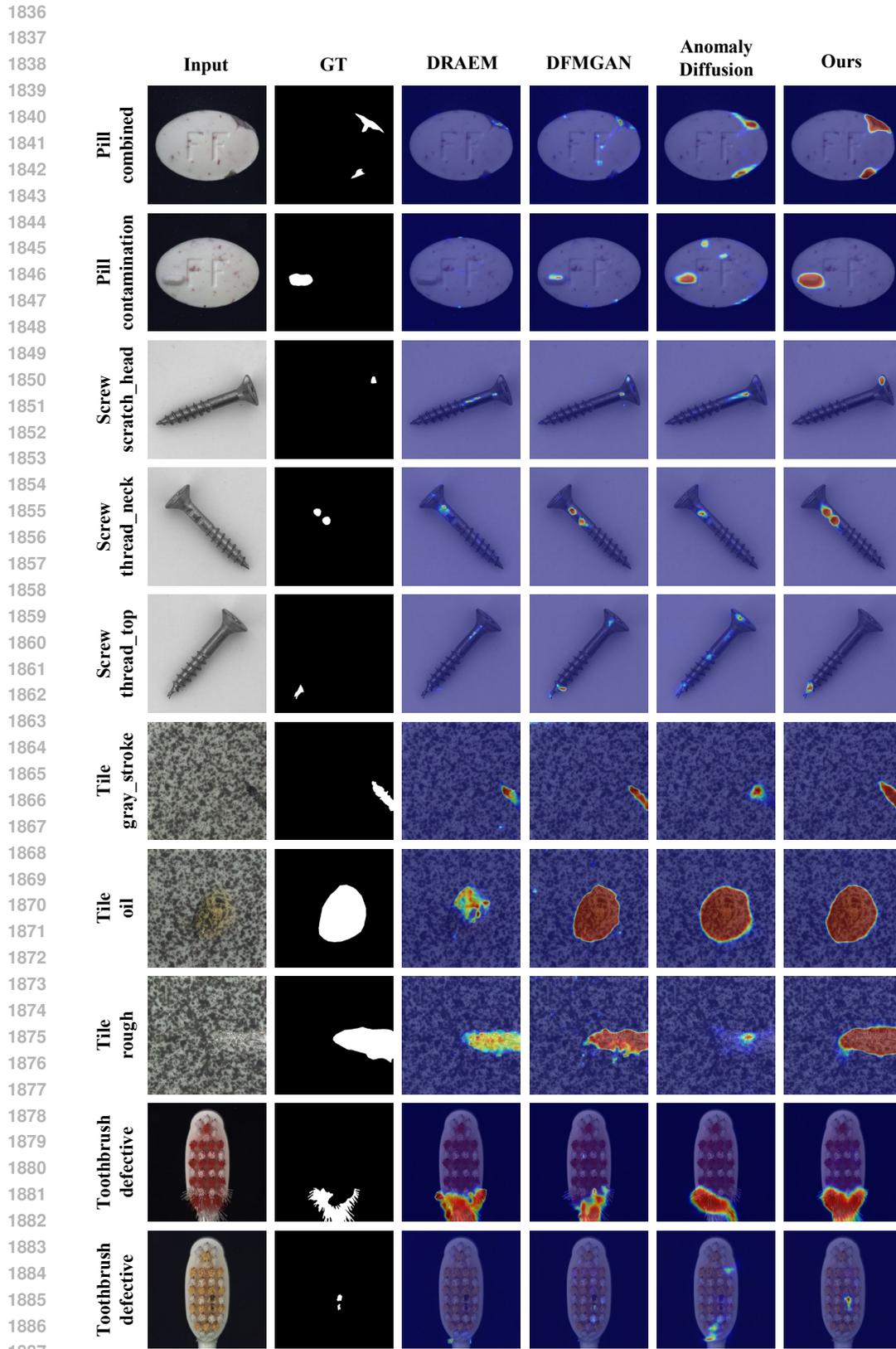


Figure 20: Comparison of the anomaly segmentation results on MVTec AD. In the figure, from top to bottom are the results for *grid*, *hazelnut*, *leather* and *metal\_nut* categories.

1834  
1835



1888  
1889  
Figure 21: Comparison of the anomaly segmentation results on MVTec AD. In the figure, from top to bottom are the results for *pill*, *screw*, *tile* and *toothbrush* categories.

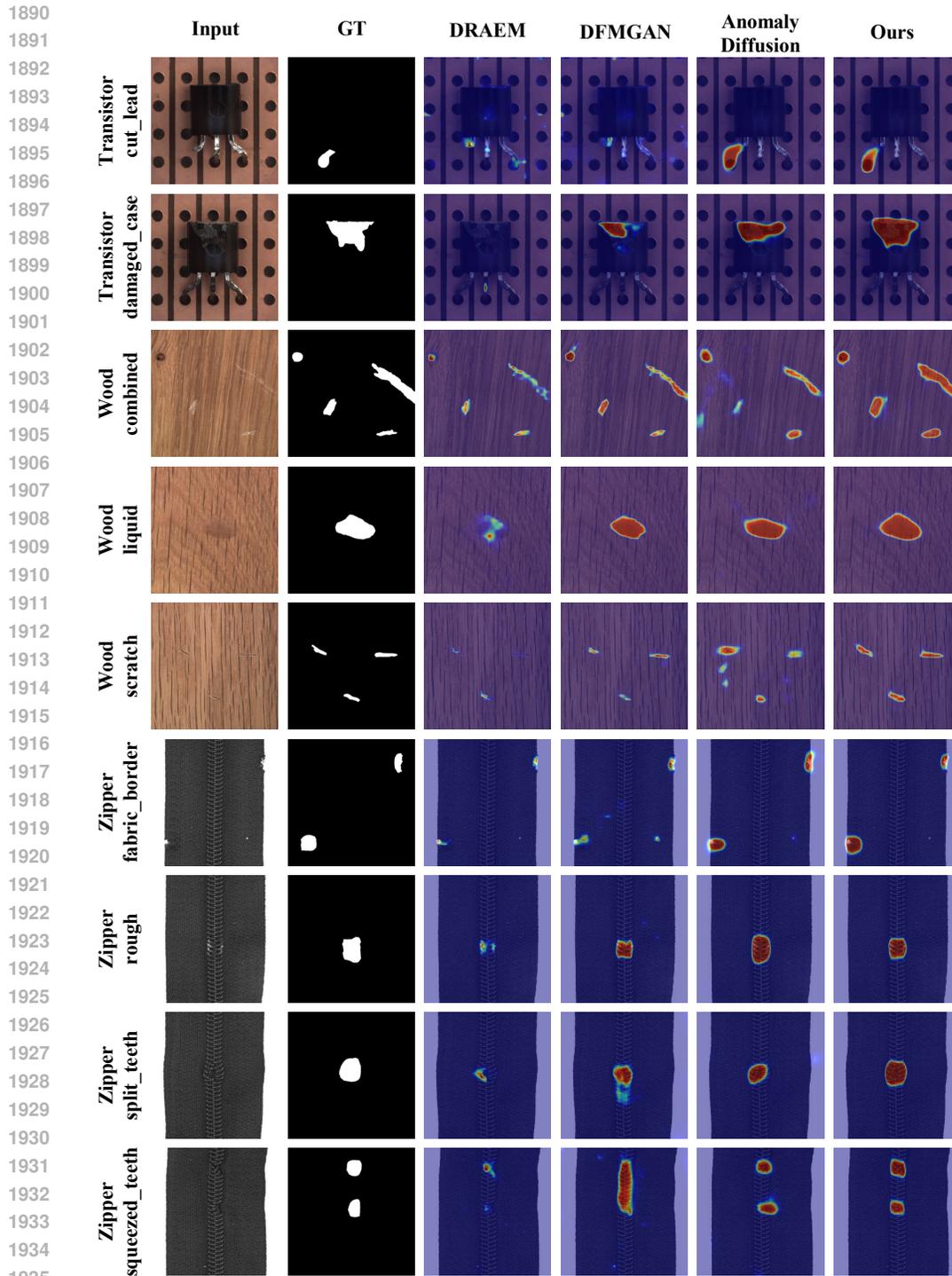


Figure 22: Comparison of the anomaly segmentation results on MVTEC AD. In the figure, from top to bottom are the results for *transistor*, *wood* and *zipper* categories.

1938  
1939  
1940  
1941  
1942  
1943

A.10 MORE QUALITATIVE COMPARISON RESULTS OF DIFFERENT SEGMENTATION MODELS TRAINED ON IMAGE-MASK PAIRS GENERATED BY SEAS

In this section, we provide further qualitative results with different segmentation models on the MVTEC AD dataset. We choose three models with different parameter quantity scopes (BiSeNet V2 (Yu et al., 2021): 3.341M, UPerNet (Xiao et al., 2018): 64.042M, LFD (Zhou et al., 2024a): 0.936M). We report the segmentation results of SeaS for varying types of anomalies in each category. Results are from Fig. 23 to Fig. 26.

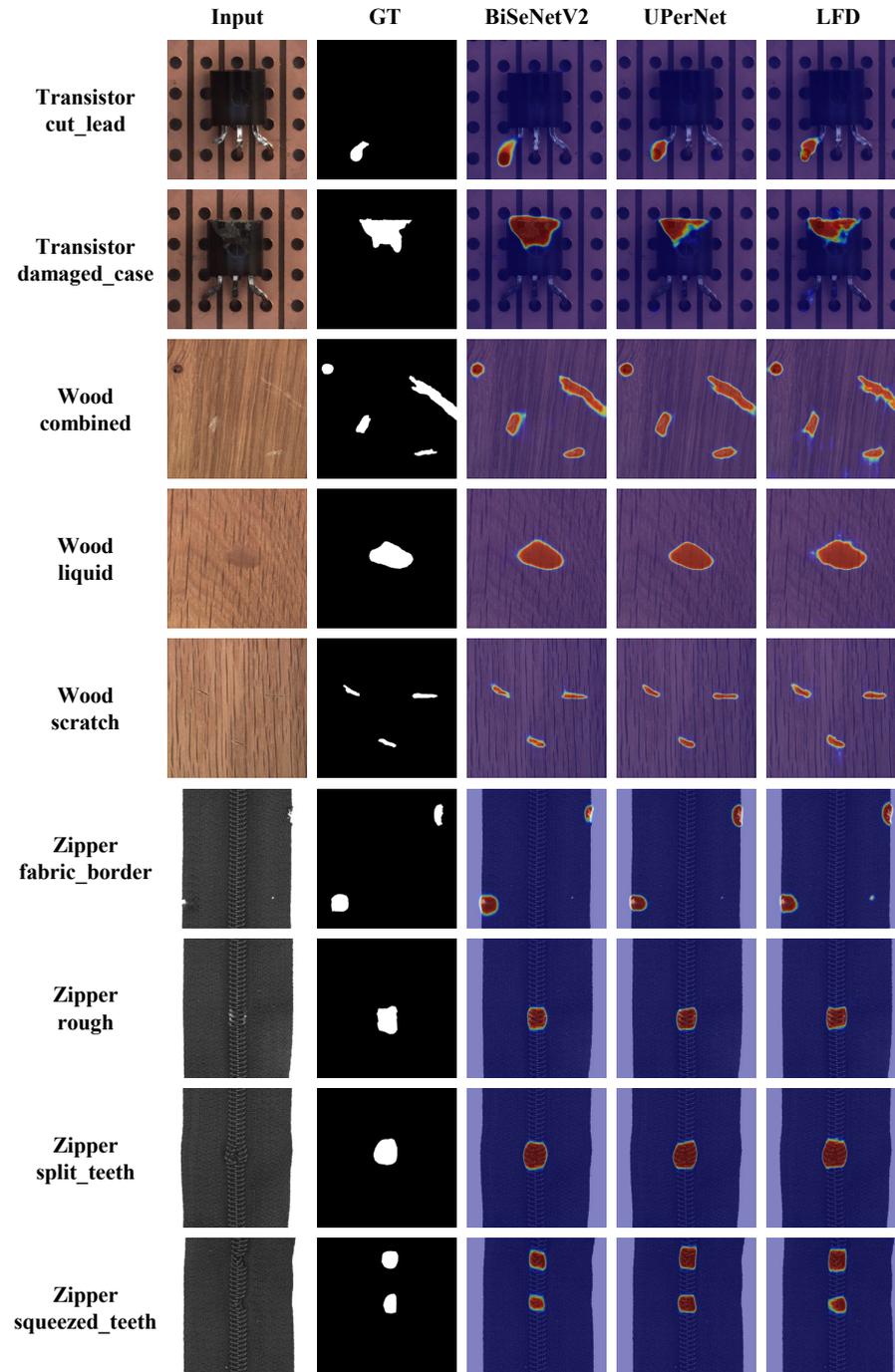


Figure 23: Qualitative comparison results with the segmentation models on MVTEC AD. In the figure, from top to bottom are the results for *transistor*, *wood*, and *zipper* categories.

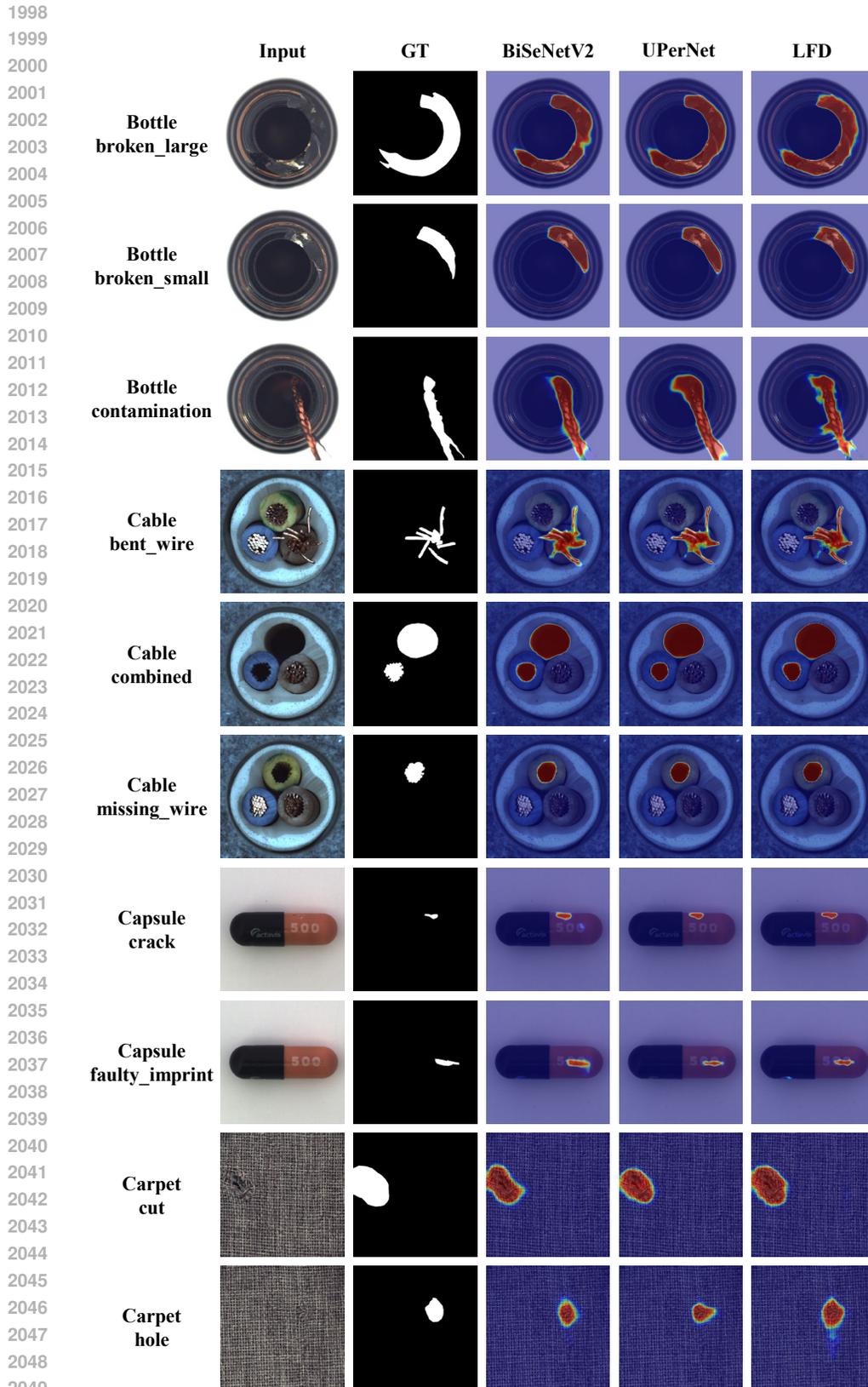


Figure 24: Qualitative comparison results with the segmentation models on MVTEC AD. In the figure, from top to bottom are the results for *bottle*, *cable*, *capsule* and *carpet* categories.

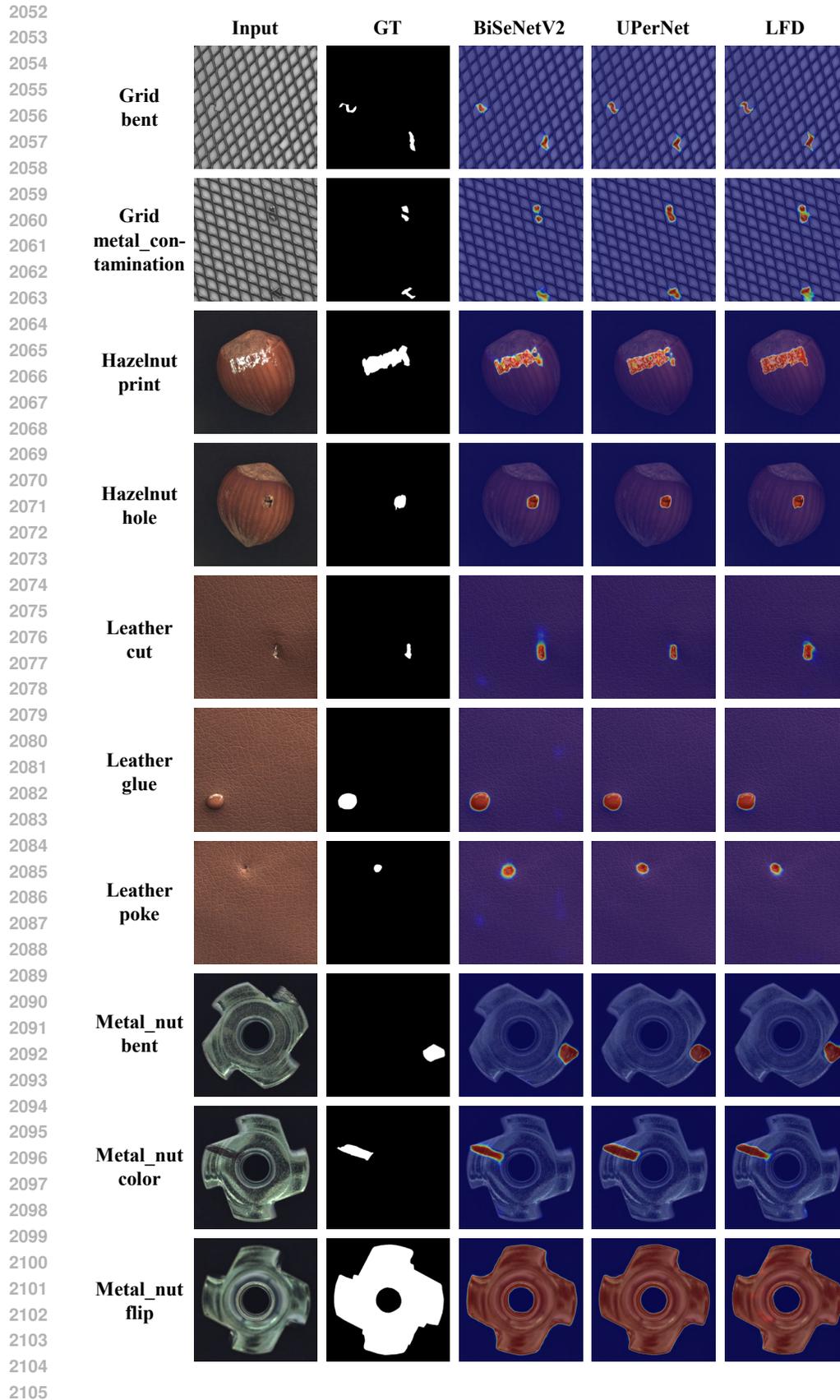


Figure 25: Qualitative comparison results with the segmentation models on MVTec AD. In the figure, from top to bottom are the results for *grid*, *hazelnut*, *leather* and *metal\_nut* categories.

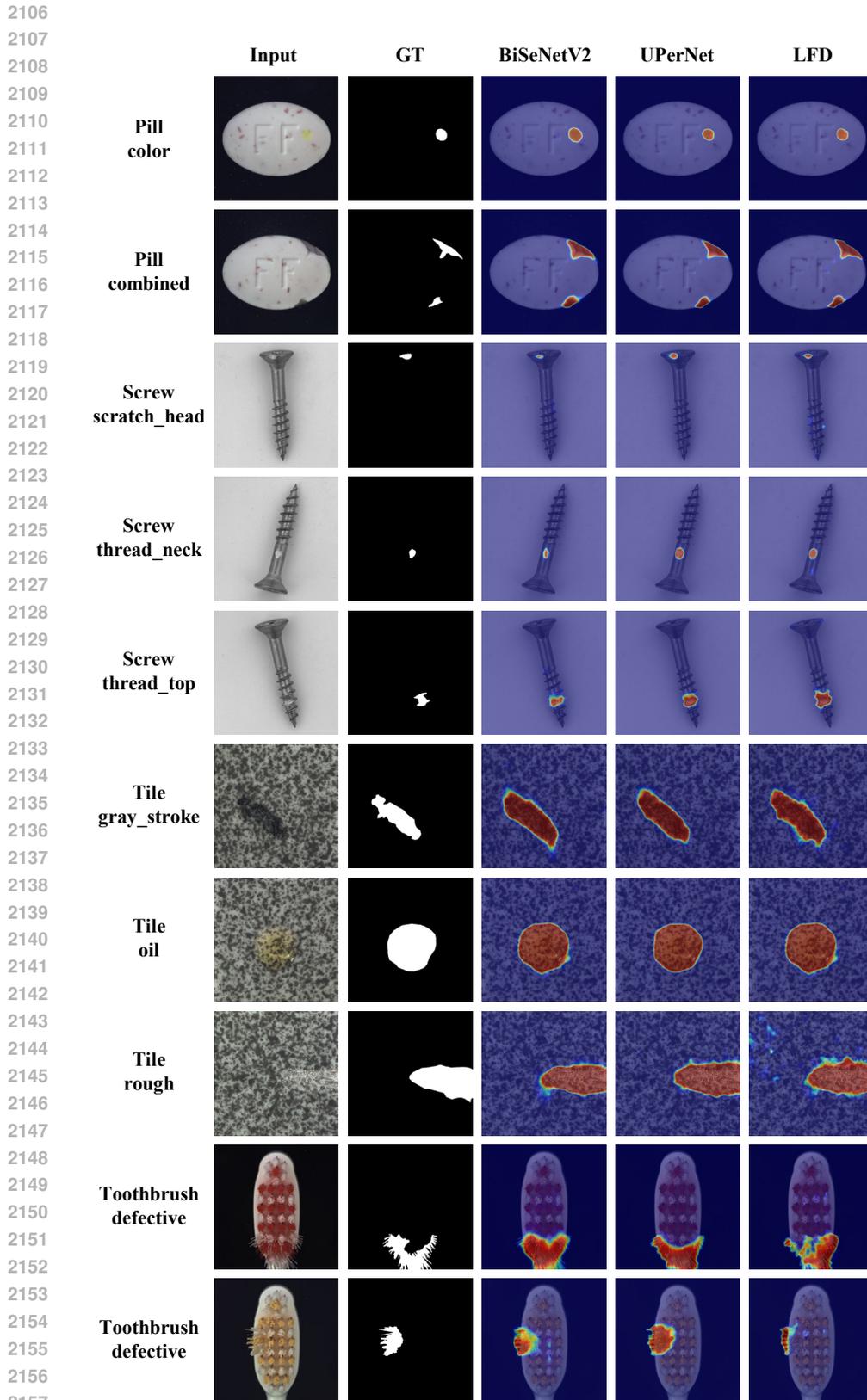


Figure 26: Qualitative comparison results with the segmentation models on MVTec AD. In the figure, from top to bottom are the results for *pill*, *screw*, *tile* and *toothbrush* categories.

A.11 ADDITIONAL VISA DATASET RESULTS

We perform experimental evaluations on the images of the VisA Dataset (Zou et al., 2022), which includes 12 product categories, each with up to 9 different anomalies.

As shown in Tab. 30 and Fig. 27, SeaS generates anomaly images with higher fidelity and diversity. Tab. 31 shows the comparisons on downstream supervised segmentation trained by the generated images. It consistently demonstrates that our method outperforms others across all the segmentation models, with an 11.71% average improvement on IoU. We report the image-level metrics in Tab. 32 and our method achieve a 5.92% gain on image-AUROC. We show the segmentation anomaly maps in Fig. 28, by using our generated image-mask pairs to train BiSeNet V2, there are fewer false positives in *chewinggum* and fewer false negatives in *pcb1* and *pipe\_fryum*.

Table 30: Comparison on IS and IC-LPIPS on VisA. Bold indicates the best performance.

Method	DFMGAN (Duan et al., 2023)		AnomalyDiffusion (Hu et al., 2024)		Ours	
	IS $\uparrow$	IC-L $\uparrow$	IS $\uparrow$	IC-L $\uparrow$	IS $\uparrow$	IC-L $\uparrow$
Average	1.25	0.25	1.26	0.25	<b>1.27</b>	<b>0.26</b>

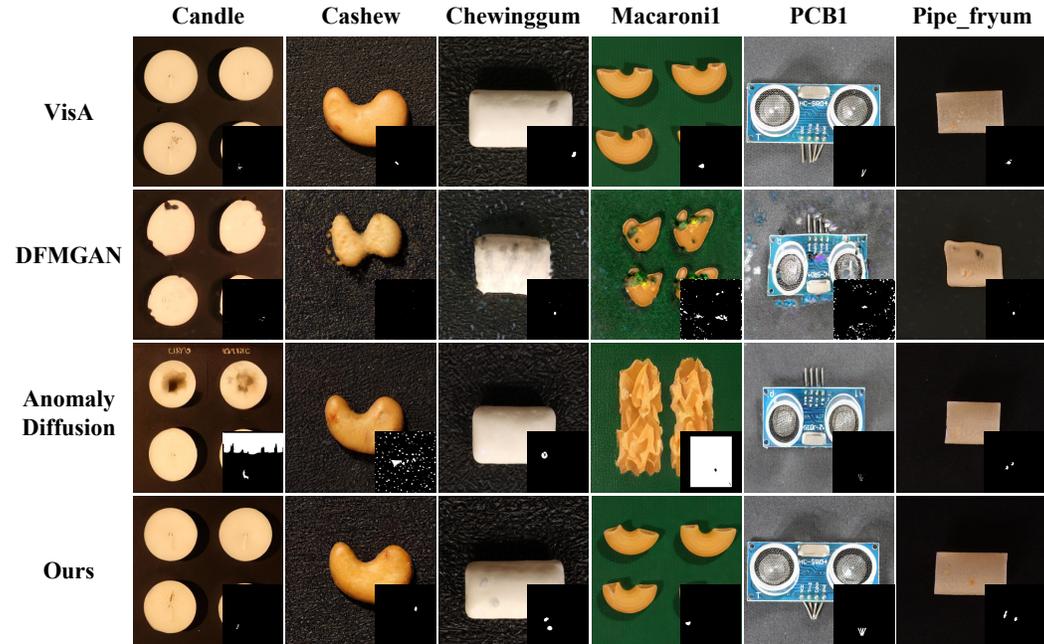


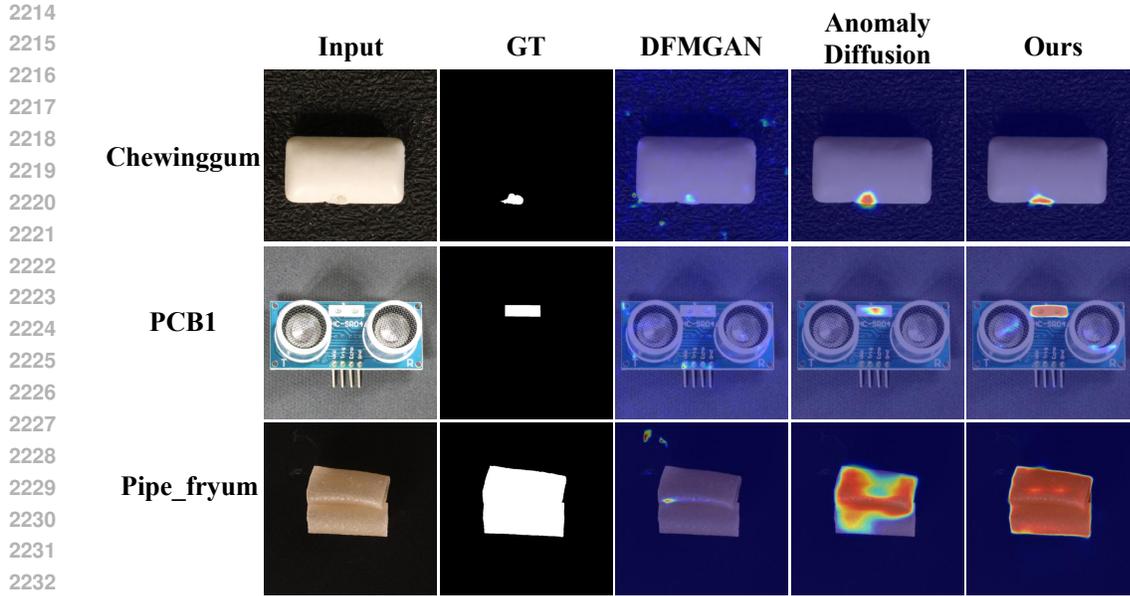
Figure 27: Visualization of the generation results on VisA. The sub-image in the lower right corner is the generated mask.

Table 31: Comparison on anomaly segmentation on VisA.

Model	DFMGAN (Duan et al., 2023)					AnomalyDiffusion (Hu et al., 2024)					Ours				
	AUROC	AP	$F_1$ -max	PRO	IoU	AUROC	AP	$F_1$ -max	PRO	IoU	AUROC	AP	$F_1$ -max	PRO	IoU
BiSeNet V2 (Yu et al., 2021)	75.91	9.17	15.00	21.49	9.66	89.29	34.16	37.93	28.09	15.93	96.03	42.80	45.41	61.29	25.93
UPerNet (Xiao et al., 2018)	75.09	12.42	18.52	27.38	15.47	95.00	39.92	45.37	44.90	20.53	97.01	55.46	55.99	58.90	35.91
LFD (Zhou et al., 2024a)	81.21	15.14	18.70	14.98	6.44	88.00	30.86	36.56	38.56	16.61	92.91	43.87	46.46	29.55	26.37
Average	77.40	12.24	17.41	21.28	10.52	90.76	34.98	39.95	37.18	17.69	<b>95.32</b>	<b>47.38</b>	<b>49.29</b>	<b>49.91</b>	<b>29.40</b>

Table 32: Comparison on image-level anomaly detection on VisA.

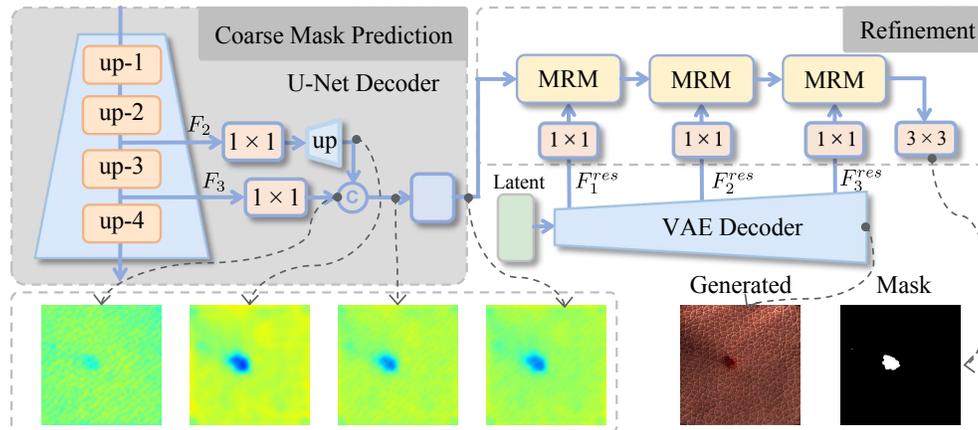
Model	DFMGAN (Duan et al., 2023)			AnomalyDiffusion (Hu et al., 2024)			Ours		
	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max
BiSeNet V2 (Yu et al., 2021)	63.07	62.63	66.48	76.11	77.74	73.13	85.61	86.64	80.49
UPerNet (Xiao et al., 2018)	71.69	71.64	70.70	83.18	84.08	78.88	90.34	90.73	84.33
LFD (Zhou et al., 2024a)	65.38	62.25	66.59	81.97	82.36	77.35	83.07	82.88	77.24
Average	66.71	65.51	67.92	80.42	81.39	76.45	<b>86.34</b>	<b>86.75</b>	<b>80.69</b>



2234 Figure 28: Qualitative anomaly segmentation results with BiSeNet V2 on VisA.

### 2235 A.12 EXPLANATION OF DISCRIMINATIVE FEATURES IN U-NET DECODER

2236  
2237  
2238 The U-Net can learn the highly discriminative features of the defect area accurately. As shown in  
2239 Fig. 29, we use the output features of the “up-2” and “up-3” layers of the decoder in U-Net, and  
2240 apply convolution blocks and concatenation operations, then we can obtain the unified coarse feature  
2241  $\hat{F} \in \mathbb{R}^{64 \times 64 \times 192}$ , which can be used to predict masks corresponding to anomaly images.



2257 Figure 29: Visualization of the U-Net decoder features in mask prediction process.

### 2258 A.13 COMPARISON WITH THE TEXTUAL INVERSION

2259  
2260  
2261 We conduct the experiment of only using the Textual Inversion (TI) (Gal et al., 2022) method to learn  
2262 the product, and the generated images are shown in Fig. 30. The TI method struggles to generate  
2263 images similar to the real product due to the limited number of learnable parameters. In contrast, for  
2264 the AIG method, the products satisfy global consistency with minor variations in local details, while  
2265 the anomalies hold randomness, so the generated products should be globally consistent with the  
2266 real products. Therefore, unlike the AG method AnomalyDiffusion (Hu et al., 2024), where the TI  
2267 method alone is sufficient to meet the anomaly generation needs, we fine-tunes the U-Net to ensure  
the global consistency of the generated products.

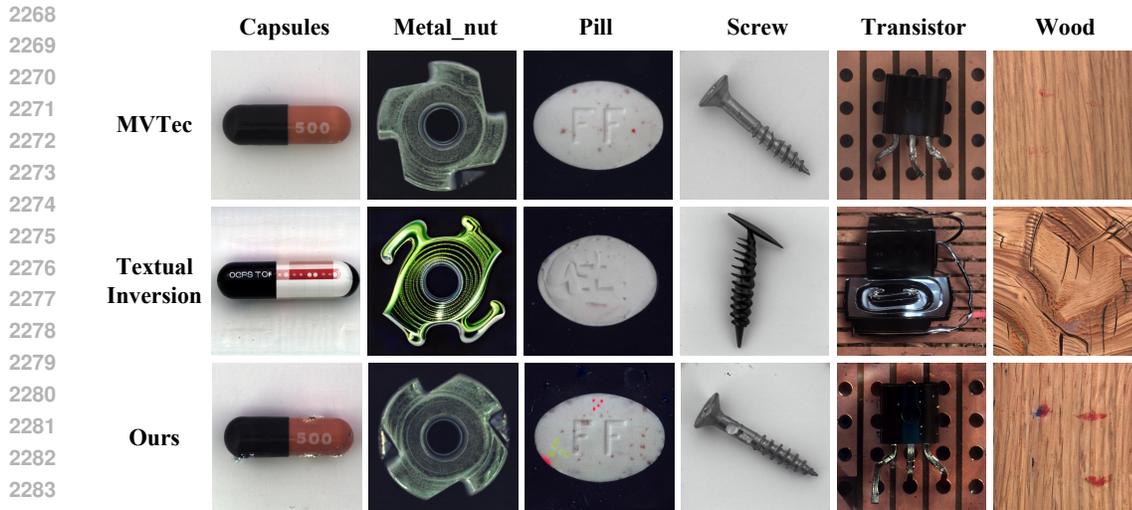


Figure 30: Qualitative comparison on the generation results with Textual Inversion.

#### A.14 MORE EXPERIMENTS ON LIGHTING CONDITIONS

We choose one defect class from peach, a product in the MVTEC3D dataset, that has significant variations in lighting conditions and backgrounds, to conduct experiments. Images with strong lighting conditions depict the top side of the peach, whereas those with weak lighting conditions show the bottom side. Consequently, the background in the images, whether the top or bottom of the peach, also differs. We selected three training sets with different lighting conditions for experiments: 1) only images from the top side with strong lighting condition, 2) only images from the bottom side with weak lighting condition, 3) half of the images from the top side with strong lighting condition, and a half from the bottom side with weak lighting condition. The generated images of different settings are shown in Fig. 31. It can be seen that SeaS is robust against lighting conditions and background variations.

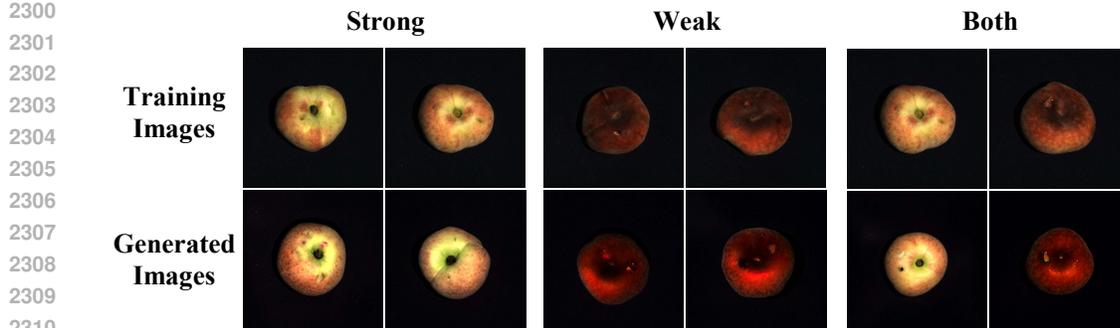


Figure 31: Visualization of the generation results on MVTEC3D AD on different lighting conditions and backgrounds. In the figure, the first row is for the training images and the second row is for the generated images.

### A.15 MORE EXPERIMENTS ON REPLACING GENERATION STRATEGIES.

We replace the abnormal generation strategy in DRAEM (Zavrtanik et al., 2021) and BGAD (Yao et al., 2023b) with the proposed generation strategy, the results are given in Tab. 33. The segmentation result demonstrates that our method outperforms the existing anomaly detection methods.

Table 33: Comparison on replacing generation strategies with anomaly detection methods on MVTEC AD.

Model	Image-level			Pixel-level			
	AUROC	AP	$F_1$ -max	AUROC	AP	$F_1$ -max	IoU
DRAEM (Zavrtanik et al., 2021)	98.00	98.45	96.34	97.90	67.89	66.04	60.30
SeaS + DRAEM	<b>99.25</b>	<b>99.66</b>	98.35	97.98	<b>77.35</b>	73.27	<b>63.99</b>
BGAD (Yao et al., 2023b)	98.31	98.05	98.27	<b>99.26</b>	73.85	77.89	60.60
SeaS + BGAD	98.44	98.18	<b>99.08</b>	<b>99.26</b>	73.85	<b>77.93</b>	60.81

### A.16 MORE VISUALIZATION RESULTS ON RECOMBINING THE DECOUPLED ATTRIBUTES FOR UNSEEN ANOMALIES.

We provide more examples in Fig. 32, where new anomalies are generated that significantly differ from the training samples in terms of color and shape. For example, we showcase *bottle\_contamination*, *hazelnut\_print*, and *tile\_gray\_stroke* with a novel shape, *wood\_color* and *metal\_nut\_scratch* with a novel color, and *pill\_crack* with a new shape, featuring multiple cracks where the training samples only exhibit a single crack. These examples demonstrate the model's ability to create unseen anomalies based on recombining the decoupled attributes.

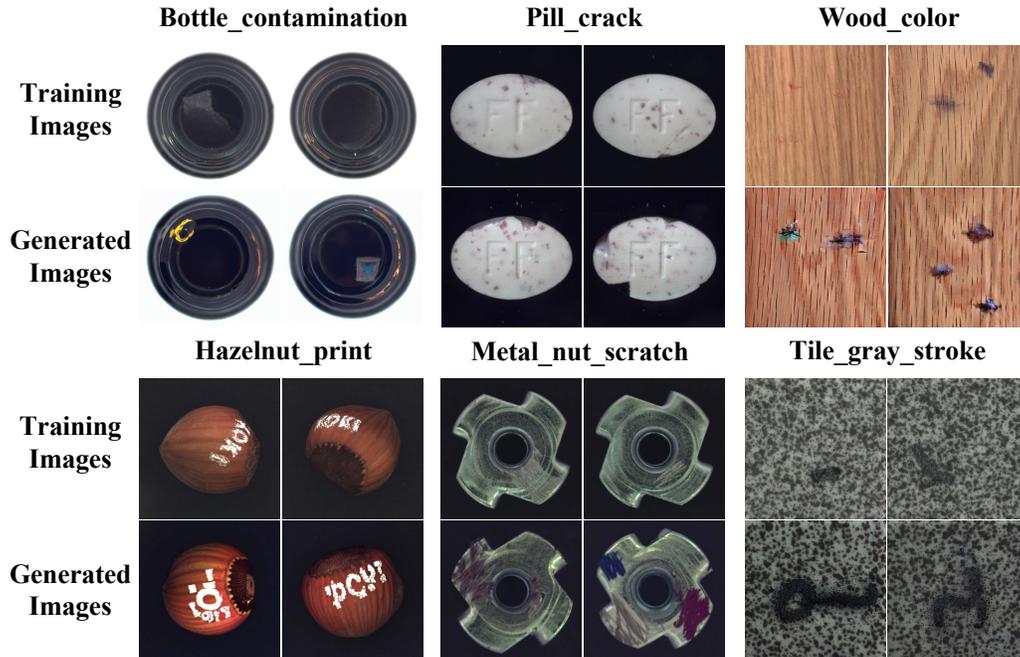


Figure 32: Visualization of the generation results for unseen anomalies on MVTEC AD.