

# A Discrepancy Aware Framework for Robust Anomaly Detection

Yuxuan Cai, Dingkang Liang, *Graduate Student Member, IEEE*, Dongliang Luo, Xinwei He, Xin Yang, *Member, IEEE*, and Xiang Bai, *Senior Member, IEEE*

**Abstract**—Defect detection is a critical research area in artificial intelligence. Recently, synthetic data-based self-supervised learning has shown great potential on this task. Although many sophisticated synthesizing strategies exist, little research has been done to investigate the robustness of models when faced with different strategies. In this paper, we focus on this issue and find that existing methods are highly sensitive to them. To alleviate this issue, we present a Discrepancy Aware Framework (DAF), which demonstrates robust performance consistently with simple and cheap strategies across different anomaly detection benchmarks. We hypothesize that the high sensitivity to synthetic data of existing self-supervised methods arises from their heavy reliance on the visual appearance of synthetic data during decoding. In contrast, our method leverages an appearance-agnostic cue to guide the decoder in identifying defects, thereby alleviating its reliance on synthetic appearance. To this end, inspired by existing knowledge distillation methods, we employ a teacher-student network, which is trained based on synthesized outliers, to compute the discrepancy map as the cue. Extensive experiments on two challenging datasets prove the robustness of our method. Under the simple synthesis strategies, it outperforms existing methods by a large margin. Furthermore, it also achieves the state-of-the-art localization performance. Code is available at: <https://github.com/caiyuxuan1120/DAF>.

**Index Terms**—Artificial intelligence, self-supervised learning, robustness

## I. INTRODUCTION

IMAGE anomaly detection plays an important role in many safety-critical areas, e.g., industrial manufacturing systems [1], [2], surveillance systems [3], and medical image analysis [4]. However, in these areas, acquiring sufficient high-quality anomaly images is generally difficult or even

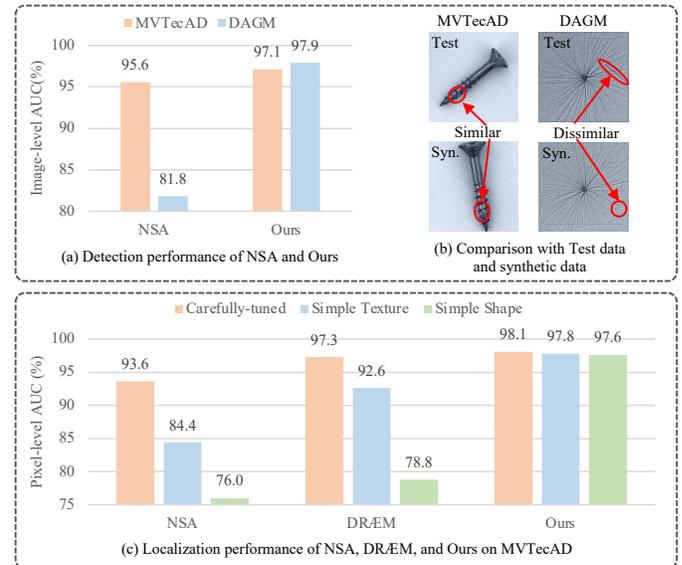


Fig. 1. (a) Detection performance comparison. (b) Real vs. Synthetic images. (c) Localization performance under diverse synthetic strategies.

impossible. This limitation hinders the effectiveness of training deep-learning models through supervised methods. As a result, it has induced growing research interest [5]–[7] in exploring approaches that focus on training anomaly detection models solely on normal images. Nevertheless, the lack of anomaly data during training poses significant challenges in extracting discriminative features for unseen anomaly data during inference.

Recently, some researchers have assumed that models trained on normal data may fail to reconstruct anomaly patterns, and have proposed identifying anomaly regions based on reconstruction failures. However, these reconstruction-based methods may be easily misled in practice because of the identical shortcut issue [8]. Another promising research direction is based on self-supervision [9]–[11]. Generally, the typical framework follows an encoder-decoder architecture (Fig. 2(a)). Normal data is first distorted to generate realistic and diverse outlier data. Subsequently, the framework is trained to differentiate normal and synthesized abnormal images. Such a strategy of exposing models to the synthesized outliers has demonstrated better empirical results, dominating current research in anomaly detection. The success of these models heavily relies on generating diverse and close-to-real anomaly images, and much effort has been endeavored to im-

Manuscript submitted 19 February 2023; revised 23 August 2023; accepted 13 September 2023. This work was supported in part by the Young Scientists Fund of the National Natural Science Foundation of China under Grant 62302188, and in part by the National Science Fund for Distinguished Young Scholars of China under Grant 62225603. Paper no. TII-23-0558. (Corresponding author: Xinwei He)

Yuxuan Cai, Dingkang Liang, and Xiang Bai are with the School of Artificial Intelligence and Automation, Huazhong University of Science and Technology, Wuhan 430074, China (e-mail: cyx\_hust@hust.edu.cn; dkliaing@hust.edu.cn; xbai@hust.edu.cn)

Dongliang Luo and Xin Yang are with the School of Electronic Information and Communications, Huazhong University of Science and Technology, Wuhan 430074, China (e-mail: ldl@hust.edu.cn; xinyang2014@hust.edu.cn)

Xinwei He is with the College of Informatics, Huazhong Agriculture University, Wuhan 430070, China (e-mail: xwhe@mail.hzau.edu.cn)

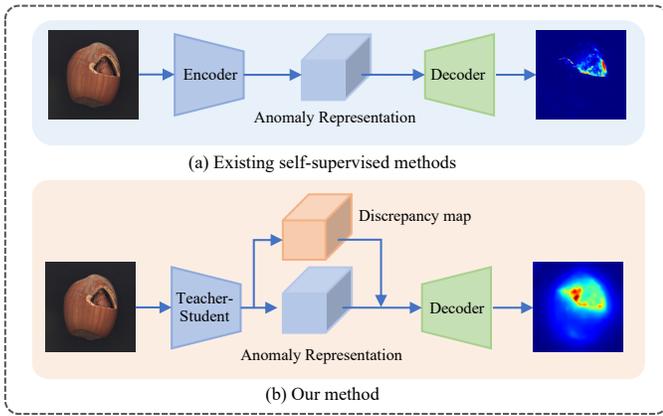


Fig. 2. Overview of existing self-supervised methods and ours.

prove anomaly synthesis strategies. For instance, DRÆM [10] uses Perlin noise to generate irregular shapes, simulating the shape of real anomalies. NSA [11] integrates Poisson image editing [12] to eliminate discontinuous borders of anomalous patterns, making the anomalies more natural.

Despite great success, few attempts have been made to study the robustness of current methods to different synthesis strategies. However, it has been noted [11] in the community that the decoder tends to overfit to the synthetic anomaly appearance during the training phase. As a result, the decision boundary tends to generalize poorly to real anomalies during inference. Moreover, the anomaly patterns generally show large variances across datasets. Therefore, the synthesis strategy customized for one dataset may not be well-suited to another. As shown in Fig. 1(b), a carefully-tuned synthesis strategy for MVTecAD [13] generates unnatural synthetic data for DAGM [14], which may help explain the dramatic performance degradation of NSA [11] on DAGM (Fig. 1(a)). Besides, our preliminary experiments on two representative self-supervision-based methods, *i.e.*, NSA [11] and DRÆM [10], also indicate their high sensitivity to different synthesis strategies (Fig. 1(c)).

To tackle the above issues, in this paper, we present a Discrepancy Aware framework (DAF) (Fig. 2(b)), which can maintain strong performance consistently across various existing anomaly synthesis techniques. The core idea behind our method is releasing the decoder from the constraint of the synthetic anomaly appearance. To accomplish this, we leverage the appearance-agnostic discrepancy map derived from a teacher-student network as guidance for the decoder. Since the discrepancy map is computed based on extracted high-level features, it is less affected by the synthetic appearance. Moreover, by processing the discrepancy directly, the decoder will focus on discriminating the normal and non-normal regions. During inference, given an anomaly image, the teacher and the student will demonstrate discrepancy in every non-normal region, while the segmentation probability map generated by the decoder will display a clear decision boundary for identifying anomalies.

To the best of our knowledge, we are the first to investigate the robustness of current frameworks to different anomaly

synthesis techniques. Compared with existing self-supervised methods, our method has the following desired properties. First, our framework incorporates the teacher-student network into the self-supervised paradigm, enhancing its capacity to produce discrepant features for anomaly regions. Second, our method encodes non-normal regions derived from the discrepancy map rather than the non-normal appearance. This approach reduces overfitting problems for the decoder during training, thus eliminating the heavy reliance on carefully-tuned anomaly synthesis techniques. As shown in Fig. 1(c), our method reaches strong performance even with a simple synthesis strategy. Besides, our method also achieves significant performance improvements over existing methods, surpassing them by a large margin in terms of localization capability on MVTecAD [13]. It also achieves the state-of-the-art detection performance on DAGM [14].

In summary, the main contributions of this paper are as follows: 1) We introduce a simple and robust self-supervised framework named **DAF** for image anomaly detection and localization, which eliminates the practical need for the complicated tuning steps for synthesis. 2) We propose to combine the teacher-student network with the self-supervised paradigm, which utilizes abundant synthesized anomaly images to learn the discrepancy features, alleviating the overfitting problem to the non-normal appearance.

## II. RELATED WORK

Early approaches [15], [16] to image anomaly detection typically work by first extracting the feature descriptors or statistical information and then calculating the anomaly scores.

Recently, with remarkable progress in deep learning, image anomaly detection based on deep learning has become a dominant direction. Below we mainly review the deep learning-based approaches, which can be roughly grouped into three categories: reconstruction-based, self-supervision-based, and knowledge distillation-based approaches.

### A. Reconstruction-based approaches

These approaches [17], [18] assume that anomalies are difficult to be reconstructed by models trained only on normal images. Thus, the anomaly regions can be spotted by examining regions with larger reconstruction errors. In these methods, autoencoders are frequently adopted as reconstruction models. For instance, Baur *et al.* [19] introduce deep spatial autoencoding architectures, which are trained by a pixel-wise reconstruction loss and an adversarial loss to improve the construction quality. Besides, generative adversarial networks [20], [21] are also attractive models for reconstructing the input image. However, the assumption behind the reconstruction methods might not always be valid, as neural networks sometimes generalize to anomalies well and yield good reconstruction results.

### B. Self-supervision-based approaches

Driven by the success of self-supervised learning in visual representation learning [22], [23], approaches under this

paradigm have emerged rapidly and advanced the state-of-the-art. For instance, [24] trains an autoencoder to reconstruct the masked image to the original one first, and then uses the reconstruction error as the anomaly score. Different from [24], which constructs masks on images, SSPCAB [25] applies the idea of masking in convolution blocks. In this way, it can be integrated into any CNN architecture.

Recent works [10], [11] prove that anomaly detection benefits from synthesized defects that are close to real ones. These works carefully design synthetic data strategies, then train segmentation models such as U-Net [26] for pixel-level prediction, ensuring that the segmentation model can learn a suitable decision boundary between normal and abnormal regions. However, the segmentation model is likely to be ineffective when there is a significant difference between the distribution of synthetic defects and actual ones. Even though DRÆM [10] tries to prevent the model from overfitting to the synthetic data by introducing anomaly-free reconstruction, as mentioned above, the reconstruction model will also fail to restore the anomalous region to normal one during inference if the anomaly manifold is unseen during training. Moreover, SPD [27] shifts its focus to self-supervised pre-training. Specifically, it proposes a novel augmentation strategy to encourage models to be locally sensitive, making the representations more suitable for the defect detection task.

### C. Knowledge distillation-based approaches

This group of approaches detects anomalies by reflecting images to different representation spaces, assuming that the representations of normal regions in different spaces are identical while those of abnormal regions will differ. The teacher-student framework is adopted to achieve the different-space reflection. For instance, prior works [5], [6], [28] train the student network to mimic the pre-trained teacher on normal images. As a result, the student and the teacher tend to hold consistency on normal regions and display discrepancies on anomalies. In the STAD [5] framework, multiple students with the same structure are trained to regress the teacher network. The discrepancies between the teacher and students, along with the variance of students, are adopted to represent the anomaly score. MKDAD [28] and STPM [6] propose distilling features from various layers of the teacher to the corresponding layers of the student. RDAD [7] utilizes reverse distillation to prohibit comparable anomaly representations in different feature spaces. SSMRKD [29] further incorporates reverse distillation with the masking strategy and constructs a two-stage framework. In the first stage, it uses reverse distillation to construct a reconstruction network, enabling the student model to accurately reconstruct normal patterns. In the second stage, the masking strategy is applied to enhance the sensitivity of the reconstruction model to anomalies, enabling effective identification and classification of abnormal patterns.

However, all the aforementioned knowledge distillation-based methods assume training models using clean data. Recently, SoftPatch [30] has addressed the scenario where the training data is contaminated with real defect data. They propose a patch-level denoising strategy to improve the robustness of the model.

TABLE I  
SYMBOL DESCRIPTION

No.	Symbol	Description
1	$T$	The Teacher Network
2	$S$	The Student Network
3	$Seg$	The Segmentation Head
4	$Aux_i$	The $i$ -th Auxiliary Head
5	$f_t^i$	Representations extracted by the teacher
6	$f_s^i$	Representations extracted by the student
7	$\bar{M}$	The Discrepancy Map
8	$M_S$	The Segmentation Probability Map
9	$M_{Score}$	The Anomaly Score Map
10	$CutP$	The synthesis strategy in CutPaste [9]
11	$DRA$	The synthesis strategy in DRÆM [10]
12	$NSA_B$	The synthesis strategy in NSA [11]

Different from previous knowledge distillation-based methods, our method incorporates the teacher-student model into a synthetic data-based self-supervised framework. It is designed to distinguish between normal patterns and abnormal ones. This objective enables our method to establish a discriminative decision boundary. Furthermore, unlike previous arts that focus on using the discrepancy for defect localization, our method leverages the teacher-student model to incorporate an appearance-agnostic cue into the self-supervised framework, thereby improving its robustness to synthetic anomaly images.

## III. OUR METHOD

### A. Overview

The framework consists of a teacher-student network, a segmentation decoder, and a series of auxiliary heads, as shown in Fig. 3. For clarity and ease of reference, Table I summarizes essential symbols and their corresponding descriptions utilized in the following text. Given an input image, the teacher-student network ( $T$ - $S$ ) is expected to demonstrate discrepancies on non-normal regions. To accomplish this, the student is trained to maintain consistency with the teacher on normal regions during training. The discrepancy maps yielded by the  $T$ - $S$  are subsequently fed into the segmentation head, along with representations of synthesized data, to localize anomalies. Finally, auxiliary heads are employed to further supervise the student model to derive more discriminative representations for anomalous regions.

### B. Teacher-Student Network

In practice, anomalies occur in various, sometimes unconstrained formats. Some can be easily discerned by their distinct textures or colors from normal patterns, while others may require contextual information to be localized. Therefore, both low-level and high-level representations are vital for accurately localizing anomalies in such challenging scenarios. In our design, the teacher-student framework follows a multi-scale knowledge distillation paradigm to represent anomalies at different levels of granularity.

Following existing knowledge distillation methods [5], [6], we adopt a powerful convolutional neural network, pre-trained on a large-scale dataset ImageNet [31], to initialize the teacher. As for the student, we choose the same architecture but

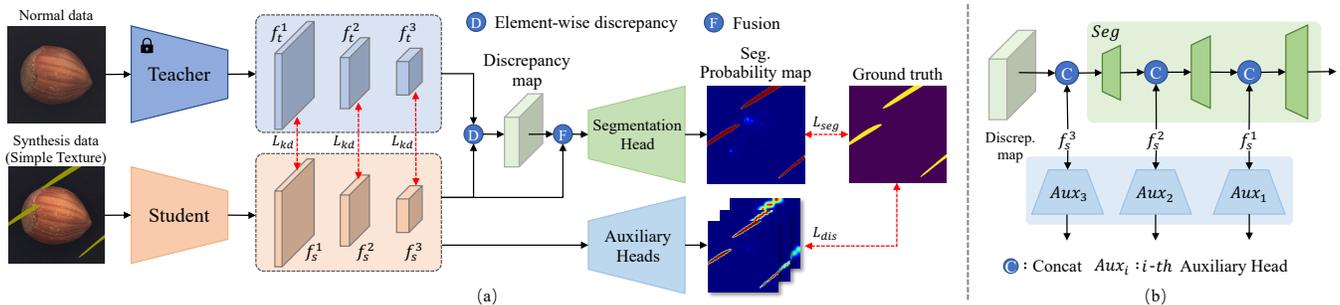


Fig. 3. (a) Overview of our method. (b) The fusion process.

initialize it randomly. Note that any off-the-shelf pre-trained networks could be adopted. Here we use ResNet18 [32], following STPM [6]. For the multi-scale knowledge distillation, we transfer the knowledge of three-stage features from the teacher to the student. The distillation process is elaborated as follows.

1) *Training strategy*: Given an anomaly-free training set  $\mathcal{D} = \{I_m\}_{m=1}^N$ , we first distort each  $I_m \in \mathcal{D}$  to obtain the synthetic anomaly image  $P_m$  using predefined anomaly strategies. The teacher network takes the anomaly-free image  $I_m$  as input, while the student network takes the corresponding synthetic anomaly image  $P_m$  as input. During training, we keep the teacher frozen and train the student to mimic the responses of the teacher on normal regions.

Following previous works [6], [28], cosine similarity is applied to measure the consistency between the representations of the teacher and the student. The cosine similarity loss is defined as:

$$L_{cos} = \sum_{i=1}^3 \left(1 - \frac{1}{N_{neg}} \sum_{\Omega_{neg}} \frac{V_t^i \cdot V_s^i}{\|V_t^i\| \cdot \|V_s^i\|}\right), \quad (1)$$

where  $i$  denotes the  $i$ -th stage.  $\Omega_{neg}$  represents the normal regions and  $N_{neg}$  indicates the total number of pixels in normal regions.  $V_t^i$  and  $V_s^i$  represent the feature vectors yielded by the teacher and the student, respectively.

However, the cosine similarity only constrains each feature vector in isolation, disregarding local contextual information. Thus we further introduce structural similarity ( $SSIM$ ) loss for compensation by considering neighboring vectors. The structural similarity is formulated as:

$$SSIM(p, q) = \frac{(2\mu_p\mu_q + \lambda_1)(2\sigma_{pq} + \lambda_2)}{(\mu_p^2 + \mu_q^2 + \lambda_1)(\sigma_p^2 + \sigma_q^2 + \lambda_2)}, \quad (2)$$

where  $\mu$  and  $\sigma$  represent the mean and variance, and  $\lambda_1$  and  $\lambda_2$  are adopted for numerical stability<sup>1</sup>. The  $SSIM$  is in the range of  $[-1, 1]$ . In particular, when  $p$  is the same as  $q$ ,  $SSIM(p, q)$  is equal to 1. Based on  $SSIM$ , the structural similarity loss is formulated as follows:

$$L_{SSIM} = \sum_{i=1}^3 (1 - SSIM(f_t^i, f_s^i)), \quad (3)$$

<sup>1</sup> $\lambda_1, \lambda_2$  are empirically set to 1e-1, 9e-4 to avoid division by zero, respectively.

where  $f_t^i$  and  $f_s^i$  represent the  $i$ -th stage feature maps of the teacher and the student, respectively. Note that we only constrain  $f_t^i$  and  $f_s^i$  on normal regions. Finally, we combine the above two losses to form the training objective of the teacher-student framework:

$$L_{kd} = L_{cos} + L_{SSIM}, \quad (4)$$

2) *The discrepancy map*: Compared with existing knowledge distillation methods [6], [7], rather than directly taking the discrepancy map as the anomaly localization result, we serve it as an additional robust cue for the segmentation task. Taken a test image  $I_m \in \mathbb{R}^{H \times W \times 3}$  as input, where  $H$  and  $W$  denote the height and width, respectively. The teacher and the student output corresponding feature representations  $f_t^i$  and  $f_s^i$ . Then the discrepancy between  $f_t^i$  and  $f_s^i$  is measured by cosine and structural similarities. Specifically, the discrepancy map at location  $(x, y)$  is defined as:

$$M_{(x,y)}^i = 2 - \frac{V_t^i \cdot V_s^i}{\|V_t^i\| \cdot \|V_s^i\|} - SSIM(f_t^i, f_s^i)_{(x,y)}, \quad (5)$$

where  $SSIM(\cdot)_{(x,y)}$  represents the structural similarity of the feature patches centered at  $(x, y)$ , and the patch size is set to  $11 \times 11$  in practice. The discrepancy map of each layer  $M^i$  is then upsampled to  $H \times W$  and summed together as the final discrepancy map  $\bar{M} \in \mathbb{R}^{H \times W}$ .

Feeding the discrepancy map into the following segmentation head brings two desirable merits. On the one hand, the segmentation head can yield more discriminative representations under this strong guidance, thereby localizing anomalies more accurately. On the other hand, the discrepancy map releases the segmentation head from being only constrained by the synthetic appearance, enhancing its perception of unseen anomalies. Such a strategy can greatly boost the robustness of our method under various simple and cheap synthesis strategies, making our method easy to use in practice.

### C. Segmentation head

The segmentation head aims at identifying the anomaly regions. As illustrated in Fig. 3(a), the segmentation head takes both the discrepancy map (*i.e.*,  $\bar{M}$ ) and anomaly features (*i.e.*,  $\{f_s^i\}_{i=1}^3$ ) as the input, where the discrepancy map indicates the location of anomalies while the anomaly features carry both low-level textual and high-level semantic information.

TABLE II  
PERFORMANCE COMPARISON ON MVTECAD

	Method	Backbone	Strategy	I-AUC	P-AUC	P-PRO	P-mAP	FLOPS	FPS	#param.
KD	RDAD*	ResNet18		97.9	97.0	92.6	54.5	4.3G	127.0	15.9M
	STPM†	ResNet18	/	95.0	96.1	86.8	47.4	3.7G	159.2	2.8M
	MKDAD*	/		86.1	88.1	75.4	23.8	5.2G	205.2	<b>0.3M</b>
	STAD	/		87.7	91.4	/	/	1948.1G	4.71	26.4M
Self-Supervision	CutPaste	/	<i>CutP</i>	96.1	96.0	/	/	/	/	/
	DRÆM*	/	<i>NSA<sub>B</sub></i>	84.1	95.1	84.0	45.9	198.4G	47.8	97.4M
	NSA*	ResNet18	<i>NSA<sub>B</sub></i>	95.6	96.3	90.5	58.7	<b>2.5G</b>	218.2	11.5M
	<b>DAF (Ours)</b>	ResNet18	<i>NSA<sub>B</sub></i>	97.1	97.5	92.4	66.0	6.8G	93.6	4.4M
	DRÆM	/	<i>DR</i>	<b>98.0</b>	97.3	/	68.4	198.4G	47.8	97.4M
	NSA*	ResNet18	<i>DR</i>	92.0	93.6	86.7	58.0	<b>2.5G</b>	218.2	11.5M
	TSDD	/	<i>DR</i>	92.8	93.9	/	60.7	/	/	/
	<b>DAF (Ours)</b>	ResNet18	<i>DR</i>	97.6	<b>98.1</b>	<b>93.0</b>	<b>68.5</b>	6.8G	93.6	4.4M

The detailed fusion process of the discrepancy map and anomaly features is shown in Fig. 3(b). The segmentation module processes them progressively in a coarse-to-fine manner via the three blocks, which is formulated as,

$$M_S = Seg_3(Seg_2(Seg_1(\overline{M} \odot f_s^3) \odot f_s^2) \odot f_s^1), \quad (6)$$

where  $Seg_1, Seg_2, Seg_3$  denote the three blocks of the segmentation module, and  $\odot$  denotes concatenation.

We apply the binary cross-entropy (BCE) loss as the segmentation loss. Note that the anomaly (positive) pixels are much fewer than normal (negative) pixels. To overcome the imbalance of positive and negative pixels, a hard negative mining strategy is adopted. Mathematically, the segmentation loss  $L_{seg}$  is defined as:

$$L_{seg} = \sum_{i \in Sub} y_i \log x_i + (1 - y_i) \log(1 - x_i), \quad (7)$$

where  $Sub$  is a subset sampled from  $M_S$ ,  $x_i$  is the predicted anomaly probability in  $Sub$ , and  $y_i$  is its corresponding label.

#### D. Auxiliary Supervision

In the teacher-student framework, the student is trained to regress representations of the teacher on normal regions. However, there is no constraint for the student on anomalous regions, which hinders the student from providing discriminative representations for the segmentation head  $Seg$ .

To enhance the discrimination capability of the student on anomalous regions, we add an auxiliary head after each student layer, as shown in Fig. 3(b). Each head  $Aux_i$  takes its corresponding anomaly feature  $f_s^i$  as input, and outputs the probability map of the corresponding size. The training strategy is the same as that of the segmentation head, meaning that we adopt the BCE loss and hard mining strategy as the discriminative loss  $L_{dis}$ . Note that auxiliary heads are discarded during inference, incurring no extra cost.

#### E. Anomaly Localization

During inference, we localize the anomalous regions from two aspects. Firstly, since the segmentation head  $Seg$  is trained to distinguish the normal distribution from that of the abnormal, the segmentation probability map is expected to localize anomalies accurately. Secondly, the discrepancy map between

the teacher and the student can also localize anomalies, since the student is likely to demonstrate inconsistency with the teacher on anomalous patterns. Since the discrepancy map is irrelevant to the synthetic appearance, combining it and the segmentation probability map desires both accurate and robust properties.

In summary, the anomalies are located by incorporating the discrepancy map  $\overline{M}$  with the segmentation probability map  $M_S$ . The final anomaly score map is formulated as:

$$M_{Score} = G(\overline{M} + \lambda M_S), \quad (8)$$

The  $G$  means Gaussian smooth.  $\lambda$  is set to 3 in practice.

## IV. EXPERIMENTS

### A. Experimental Setup

1) *Datasets*: MVTECAD [13] contains 5,354 images, including ten object categories and five texture categories. The training set has 3,629 normal images, while the testing set contains 1,725 images, covering both normal and anomaly images. Pixel annotations are provided for the anomaly areas.

DAGM [14] includes ten categories of texture images. The training set contains anomaly images, and weak annotations are provided for anomalous regions. During training, we only use normal images. Since the annotations are coarse, we do not evaluate the localization performance on DAGM.

2) *Model Training*: All the images are resized to  $256 \times 256$ . The weights of the teacher are frozen, and the remaining components are trained using AdamW for 1,200 epochs. The batch size is set to 8, and the weight decay is set to  $1e-5$ . The learning rate is gradually increased to  $2e-4$  in 50 epochs and multiplied by 0.2 after 700 and 1,000 epochs. We first follow the synthetic strategies in DRÆM [10] and NSA [11] and then adopt a series of simple synthetic approaches to further investigate the effectiveness of our method. Note that DRÆM [10] introduces an external dataset (*i.e.* DTD [33]) for synthesizing.

3) *Evaluation Metrics*: Following Salehi *et al.* [28], we evaluate the performance of anomaly detection and localization by image-level AUC (I-AUC) and pixel-level AUC (P-AUC), respectively. Meanwhile, following [13] and [10], we also focus on the Per-Region-Overlap (P-PRO) and mean Average Precision (P-mAP) to further evaluate localization accuracy.

TABLE III  
LOCALIZATION PERFORMANCE COMPARISON UNDER SIMPLE AND CHEAP SYNTHESIS STRATEGIES

	Metric	Method	Carp.	Grid	Leath.	Tile	Wood	Bottle	Cable	Caps.	Haze.	Metal.	Pill	Screw	Tooth.	Trans.	Zip.	Mean	
Simple Texture	P-AUC	DRÆM	/	/	/	/	/	/	/	/	/	/	/	/	/	/	/	92.6	
		NSA*	91.3	80.3	90.6	95.0	83.2	92.9	76.4	83.8	91.4	96.7	80.9	73.5	96.8	68.6	64.5	84.4	
		DAF	<b>98.8</b>	<b>98.8</b>	<b>99.6</b>	<b>97.7</b>	<b>96.8</b>	<b>98.7</b>	<b>97.2</b>	<b>97.4</b>	<b>99.0</b>	<b>99.1</b>	<b>98.0</b>	<b>99.0</b>	<b>99.1</b>	<b>89.8</b>	<b>98.3</b>	<b>97.8</b>	
	P-PRO	DRÆM	/	/	/	/	/	/	/	/	/	/	/	/	/	/	/	/	/
		NSA*	79.9	54.2	88.7	<b>93.9</b>	79.8	78.1	62.9	66.8	89.5	91.1	78.7	38.1	75.3	53.7	36.4	71.1	
		DAF	<b>94.2</b>	<b>94.2</b>	<b>98.5</b>	92.7	<b>92.9</b>	<b>94.5</b>	<b>89.9</b>	<b>84.3</b>	<b>95.6</b>	93.9	<b>92.4</b>	<b>95.1</b>	90.3	<b>78.7</b>	<b>91.7</b>	<b>91.9</b>	
P-mAP	DRÆM	/	/	/	/	/	/	/	/	/	/	/	/	/	/	/	/	56.5	
	NSA*	47.7	8.2	42.4	<b>89.1</b>	68.0	62.3	16.8	30.1	63.9	91.5	28.3	2.2	55.7	22.9	13.8	42.9		
	DAF	<b>67.1</b>	<b>43.1</b>	<b>66.8</b>	84.3	<b>73.2</b>	<b>81.6</b>	<b>55.8</b>	<b>41.0</b>	<b>67.2</b>	<b>93.8</b>	<b>54.1</b>	<b>34.2</b>	<b>66.0</b>	<b>55.3</b>	<b>49.7</b>	<b>62.2</b>		
Simple Shape	P-AUC	DRÆM*	79.8	97.6	91.5	80.2	74.6	74.5	71.3	63.4	81.3	73.0	72.2	81.7	87.5	71.1	82.4	78.8	
		NSA*	76.6	95.4	95.3	77.4	78.0	73.6	61.4	63.4	87.2	59.3	67.6	82.8	79.5	67.4	75.8	76.0	
		DAF	<b>98.6</b>	<b>99.3</b>	<b>99.7</b>	<b>98.6</b>	<b>96.4</b>	<b>98.5</b>	<b>96.7</b>	<b>95.7</b>	<b>98.8</b>	<b>98.5</b>	<b>97.5</b>	<b>98.8</b>	<b>99.0</b>	<b>88.1</b>	<b>99.1</b>	<b>97.6</b>	
	P-PRO	DRÆM*	70.8	94.9	89.4	59.4	71.0	61.3	30.6	55.6	78.1	56.3	75.8	66.9	78.9	53.1	70.9	67.5	
		NSA*	65.1	88.0	93.2	53.3	78.0	61.1	38.5	60.8	87.9	25.5	53.0	64.5	71.6	49.7	59.9	63.3	
		DAF	<b>96.1</b>	<b>97.3</b>	<b>99.0</b>	<b>95.0</b>	<b>93.1</b>	<b>95.1</b>	<b>90.3</b>	<b>81.7</b>	<b>96.5</b>	<b>96.0</b>	<b>93.3</b>	<b>94.2</b>	<b>93.4</b>	<b>75.2</b>	<b>96.0</b>	<b>92.8</b>	
P-mAP	DRÆM*	23.5	49.1	55.5	29.4	32.5	36.8	5.8	6.3	31.2	22.1	18.8	6.9	25.5	11.3	33.3	25.9		
	NSA*	12.8	45.3	34.2	19.1	53.8	43.4	4.7	13.3	42.5	13.7	6.2	9.6	16.9	12.3	30.1	23.9		
	DAF	<b>70.1</b>	<b>57.4</b>	<b>69.8</b>	<b>88.2</b>	<b>65.9</b>	<b>80.1</b>	<b>45.5</b>	<b>24.1</b>	<b>60.9</b>	<b>86.8</b>	<b>72.1</b>	<b>42.1</b>	<b>54.5</b>	<b>51.7</b>	<b>71.1</b>	<b>62.7</b>		
Simple Texture-Shape	P-AUC	DRÆM*	58.6	95.0	71.6	85.4	69.2	80.2	47.2	58.4	55.1	66.5	46.9	80.8	59.8	57.9	59.9	66.2	
		NSA*	59.7	60.6	67.5	55.7	58.4	51.4	52.2	50.0	59.8	52.7	48.0	58.1	53.5	47.8	60.2	55.7	
		DAF	<b>99.0</b>	<b>98.9</b>	<b>99.6</b>	<b>98.0</b>	<b>96.2</b>	<b>98.4</b>	<b>96.2</b>	<b>97.6</b>	<b>98.6</b>	<b>98.7</b>	<b>97.5</b>	<b>99.1</b>	<b>98.9</b>	<b>88.2</b>	<b>98.6</b>	<b>97.6</b>	
	P-PRO	DRÆM*	40.0	90.3	69.2	58.2	44.2	57.6	30.6	43.4	38.4	60.1	53.9	64.3	70.7	43.7	38.9	53.6	
		NSA*	35.2	32.4	56.7	36.0	53.6	24.7	18.6	23.9	23.5	22.6	32.5	25.0	20.9	22.5	29.3	30.5	
		DAF	<b>96.1</b>	<b>95.2</b>	<b>98.8</b>	<b>93.4</b>	<b>93.1</b>	<b>94.5</b>	<b>88.9</b>	<b>87.6</b>	<b>94.1</b>	<b>92.8</b>	<b>94.4</b>	<b>95.5</b>	<b>91.9</b>	<b>75.3</b>	<b>92.9</b>	<b>92.3</b>	
P-mAP	DRÆM*	6.5	<b>47.1</b>	42.7	46.4	15.6	42.0	4.2	8.3	19.9	31.2	14.3	13.8	25.2	10.8	17.3	23.0		
	NSA*	2.9	3.2	8.2	14.3	16.8	7.2	3.0	1.0	2.7	12.7	3.3	0.4	1.9	5.0	3.9	5.8		
	DAF	<b>71.6</b>	42.5	<b>67.3</b>	<b>85.1</b>	<b>63.1</b>	<b>81.6</b>	<b>44.7</b>	<b>35.7</b>	<b>60.2</b>	<b>88.4</b>	<b>72.7</b>	<b>29.9</b>	<b>53.3</b>	<b>53.7</b>	<b>53.0</b>	<b>60.2</b>		

4) *Baselines*: We mainly compare DAF with knowledge distillation-based methods and self-supervision-based methods. For ease of description, we abbreviate these two methods in the following tables as KD and Self-Sup, respectively.

(1) Knowledge distillation-based methods: STAD [5], STPM [6], MKDAD [28], RDAD [7]. These methods compare activations of the teacher and student, where features of anomaly regions are distinct.

(2) Self-supervision-based methods: CutPaste [9], DRÆM [10], NSA [11]. These methods introduce extra supervision through synthetic anomalies to help proxy tasks such as segmentation.

### B. Evaluation on MVTecAD using complicated strategies

1) *Quantitative analysis*: Table II reports the detection and localization performance on the MVTecAD [13] dataset<sup>2</sup>. In our method, the image-level score is acquired by the mean of the top 50 values of the anomaly score map  $M_{Score}$ . Compared with knowledge distillation methods, our method achieves a comparable result of 97.6% I-AUC with RDAD [7] on the detection task. In terms of the localization task, benefiting from the discriminative capacity of the segmentation head  $Seg$ , our approach surpasses RDAD by 1.1% P-AUC and 14.0% P-mAP. Compared with self-supervised methods, while using the synthesis strategy of  $NSA_B$ , we exceed NSA [11] by 1.5% I-AUC and 7.3% P-mAP, respectively. Furthermore, compared with DRÆM [10], our method achieves significant improvements of 13.0% I-AUC and 20.1% P-mAP. Notably, when using the synthesis strategy of  $DRA$ , our detection

<sup>2</sup>\* represents the reproduced result using the official code. † denotes the reproduced results using the unofficial code.

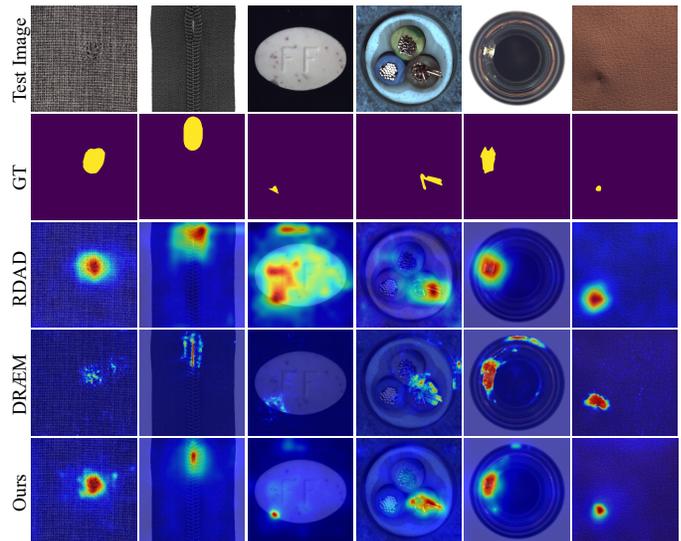


Fig. 4. Visualization of localization results.

performance is slightly lower than DRÆM [10] by 0.4%. However, our localization performance exceeds it by 0.8% P-AUC. Please note that the computation complexity (FLOPS), the model size (#param), and the FPS of DRÆM are much higher than ours. When training with synthetic data, the supervised method TSDD [34] obtains relatively poor performance, indicating that the distribution gap with real anomalies may hurt its generalization.

2) *Qualitative analysis*: Fig. 4 shows localization results of the existing state-of-the-art knowledge distillation method RDAD [7], the self-supervised method DRÆM [10] and our

TABLE IV  
DETECTION PERFORMANCE ON DAGM [14] UNDER DIFFERENT SYNTHESIS STRATEGIES

	Method	Strategy	Class1	Class2	Class3	Class4	Class5	Class6	Class7	Class8	Class9	Class10	Mean
KD	RDAD*	/	95.3	99.2	80.6	100.0	78.1	91.6	99.5	63.1	95.3	99.1	90.2
	STPM†	/	86.1	98.7	88.7	100.0	82.3	94.4	99.8	65.6	91.7	99.7	90.7
Self-Supervision	DRÆM	/	/	/	/	/	/	/	/	/	/	/	99.0
	NSA*	<i>DRA</i>	90.6	99.2	99.8	<b>100.0</b>	<b>98.1</b>	99.9	<b>100.0</b>	<b>99.5</b>	51.4	99.7	93.8
	DAF		<b>99.5</b>	<b>99.8</b>	<b>100.0</b>	<b>100.0</b>	96.7	<b>100.0</b>	<b>100.0</b>	<b>97.1</b>	<b>98.9</b>	<b>99.9</b>	<b>99.2</b>
	DRÆM*	<i>NSA<sub>B</sub></i>	45.1	<b>100.0</b>	88.5	80.0	84.1	<b>100.0</b>	<b>100.0</b>	83.9	91.2	87.9	86.1
	NSA*		52.6	99.3	85.1	97.8	86.7	<b>100.0</b>	57.4	83.2	62.1	93.8	81.8
	DAF		<b>85.7</b>	<b>100.0</b>	<b>99.6</b>	<b>99.8</b>	<b>93.8</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>99.7</b>	<b>99.9</b>	<b>99.9</b>
	DRÆM*	<i>Simple Texture</i>	54.9	45.6	98.8	68.8	<b>99.3</b>	57.0	93.4	87.4	53.4	72.6	73.1
	NSA*		50.6	50.8	80.5	51.4	72.7	64.6	<b>100.0</b>	<b>97.1</b>	53.3	76.1	69.7
	DAF		<b>97.6</b>	<b>97.1</b>	<b>99.7</b>	<b>100.0</b>	76.1	<b>96.7</b>	<b>100.0</b>	88	<b>91.8</b>	<b>99.7</b>	<b>94.7</b>
	DRÆM*	<i>Simple Shape</i>	<b>99.4</b>	90.2	89.1	96.7	<b>100.0</b>	86.9	<b>100.0</b>	92.3	48.2	77.7	88.1
	NSA*		52.2	44.2	36.6	29.6	28.5	58.8	20.5	46.5	52.8	25.1	39.5
	DAF		93.6	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	97.8	<b>100.0</b>	<b>100.0</b>	<b>97.0</b>	<b>94.8</b>	<b>99.9</b>	<b>98.3</b>

TABLE V  
ABLATION STUDY ON THE EFFECTIVENESS OF EACH COMPONENT

Components	I-AUC	P-AUC	P-PRO	P-mAP
<i>only T-S</i>	93.6	96.3	86.1	48.6
<i>only Seg</i>	95.1	95.2	75.1	62.5
<i>W/O Aux</i>	<b>97.7</b>	98.0	92.6	66.6
DAF	97.6	<b>98.1</b>	<b>93.0</b>	<b>68.5</b>

TABLE VI  
ABLATION STUDY ON THE DISCREPANCY MAP

Strategy	$\bar{M}$	$M_S$	I-AUC	P-AUC	P-PRO	P-mAP
<i>DRA</i>	✓		95.4	97.4	90.7	52.8
		✓	96.8	95.5	78.4	67.1
<i>Simple Texture</i>	✓	✓	<b>97.6</b>	<b>98.1</b>	<b>93.0</b>	<b>68.5</b>
			97.4	97.4	90.8	52.1
<i>Simple Shape</i>	✓	✓	95.3	89.9	55.6	52.9
			<b>97.5</b>	<b>97.8</b>	<b>91.9</b>	<b>62.2</b>
<i>Simple Shape</i>	✓		92.9	97.0	89.9	48.7
		✓	94.0	80.8	72.2	40.0
	✓	✓	<b>97.0</b>	<b>97.6</b>	<b>92.8</b>	<b>62.7</b>

method. Both DRÆM and our method are trained using the *DRA* strategy. It can be observed that RDAD can roughly localize anomalies, but the localization areas tend to be larger than the ground truth (Row 3, Columns 2,4,5,6). DRÆM can sometimes deliver clear decisions for normal and abnormal areas (Row 4, Column 5), but it also fails in some cases (Row 4, Columns 1-4). The last row suggests that our method achieves the most accurate localization, performing well across various anomalies and closely matching the ground truth.

### C. Evaluation on MVTecAD using simple strategies

*Simple Texture* means that the textures of the synthetic anomalies are random colors, rather than natural textures sampled from the DTD dataset [33]. *Simple Shape* indicates that the shapes of synthetic anomalies are rectangles, in contrast to irregular shapes. *Simple Texture-Shape* replaces shapes and textures with rectangles and randomly sampled colors.

1) *Quantitative analysis*: Table III reports the results. Under the *Simple Texture* strategy, it can be observed that our method achieves competitive results of 97.8% P-AUC, 91.9% P-PRO

TABLE VII  
IMPACT OF THE DISCREPANCY MAP  $\bar{M}$  ON SEGMENTATION

Strategy	$\bar{M}$	I-AUC	P-AUC	P-PRO	P-mAP
<i>DRA</i>	✓	<b>97.6</b>	98.0	92.9	67.3
		<b>97.6</b>	<b>98.1</b>	<b>93.0</b>	<b>68.5</b>
<i>Simple Texture</i>		96.5	97.7	91.3	60.1
	✓	<b>97.5</b>	<b>97.8</b>	<b>91.9</b>	<b>62.2</b>
<i>Simple Shape</i>		96.5	97.5	92.6	62.2
	✓	<b>97.0</b>	<b>97.6</b>	<b>92.8</b>	<b>62.7</b>

and 62.2% P-mAP, while DRÆM [10] and NSA [11] perform relatively lower. Specifically, compared with DRÆM, we bring 5.2% (97.8% vs. 92.6%) P-AUC and 5.7% (62.2% vs. 56.5%) P-mAP gains, and outperform NSA by 13.4% P-AUC and 19.3% P-mAP. Similarly, compared with these two prior arts, we also achieve significant improvements under the *Simple Texture* and *Simple Texture-Shape* strategies. These results demonstrate the superior robustness of our framework.

2) *Discussions*: NSA [11] carefully designs a synthesis strategy to simulate natural images. Since the optimization process is only influenced by the anomalous appearance, the model is likely to overfit to the synthetic anomalies. Hence, it can perform well if the real anomalies are similar to the synthetic ones. However, it is challenging to localize anomalies when there is a distribution discrepancy between them. DRÆM [10] attempts to solve the overfitting issue by jointly using the synthetic anomaly data and anomaly-free reconstruction. However, the reconstruction model also suffers from the overfitting problem, failing to restore anomalies when unseen anomalies occur. Our method benefits from the discrepancy map, which inherently indicates the location of anomaly regions, providing strong evidence for the segmentation head. This alleviates its dependence on the synthetic appearance. Moreover, the supplement of the discrepancy map can further promote robustness.

### D. Evaluation on DAGM

Table IV shows the detection performance on DAGM [14]. The two knowledge distillation methods RDAD [7] and STPM [6] achieve comparable performance (90.2% vs.

TABLE VIII

ABLATION STUDY ON THE HYPERPARAMETER UNDER *DRA* STRATEGY

Method	DTD	Perlin	$\beta$	I-AUC	P-AUC	P-PRO	P-mAP
DRÆM	✓	✓	✓	98.0	97.3	/	68.4
	✓	✓		97.4	95.0	/	64.5
DAF	✓	✓	✓	97.6	98.1	93.0	68.5
	✓	✓		97.6	97.8	92.4	64.5

TABLE IX

COMPARISON WITH THE MODEL ENSEMBLE

Strategy	Method	I-AUC	P-AUC	P-PRO	P-mAP
<i>DRA</i>	<i>Ensem.</i>	96.8	97.7	91.7	65.1
	DAF	<b>97.6</b>	<b>98.1</b>	<b>93.0</b>	<b>68.5</b>
<i>Simple Texture</i>	<i>Ensem.</i>	95.8	97.4	90.1	58.4
	DAF	<b>97.5</b>	<b>97.8</b>	<b>91.9</b>	<b>62.2</b>
<i>Simple Shape</i>	<i>Ensem.</i>	93.5	97.5	91.7	61.2
	DAF	<b>97.0</b>	<b>97.6</b>	<b>92.8</b>	<b>62.7</b>

90.7%). When training with *DRA*, the self-supervised methods DRÆM [10] and NSA [11] outperform knowledge distillation methods by a large margin. Our method obtains the state-of-the-art performance of 99.2%. We assume that *DRA* can synthesize anomalies that match real ones during training, helping the self-supervised methods to yield clear decision boundaries to classify the normal and abnormal regions. However, these methods demonstrate unsatisfactory results when the synthetic strategies are simpler, e.g., *Simple Texture* and *Simple Shape*. In contrast, our method demonstrates strong robustness to these synthesis strategies. The improvement is thanks to the discrepancy map, which reveals the location of anomalies, releasing the segmentation task from the constraint of the anomaly appearance. Additionally, the supplement of the discrepancy map during inference has further enhanced robustness.

### E. Ablation Study

For simplicity and fairness, we conduct all the ablation studies on MVTECAD [13].

1) *Effectiveness of different components*: Table V displays the impacts of different components. *Only T-S* indicates that only the teacher-student framework remains, and only the discrepancy map is employed for evaluation. We also only keep the student and the segmentation head to build a segmentation network (*Only Seg*). To demonstrate the effectiveness of the auxiliary heads, we remove them from the pipeline (*W/O Aux*). Compared with *only Seg*, *only T-S* shows its advantage in P-AUC and P-PRO. In contrast, *only Seg* brings higher P-mAP. *W/O Aux* brings gains in both detection and localization. By incorporating the auxiliary heads during training (Ours), we achieve the best localization performance. This observation indicates that the auxiliary heads promote the student to yield more discriminative representations for anomalies.

2) *Effectiveness of the discrepancy map*: We further conduct several settings to prove that the discrepancy map  $\bar{M}$  is a useful supplement to the segmentation probability map  $M_S$ . Table VI shows that with  $\bar{M}$ , we have observed 2.6% P-AUC, 14.6% P-PRO and 1.4% P-mAP improvements under

TABLE X

IMPACT OF INITIALIZATION FOR THE STUDENT

Strategy	Initial Param.	I-AUC	P-AUC	P-PRO	P-mAP
<i>DRA</i>	Pretrained	<b>97.6</b>	97.9	91.9	64.3
	Random	<b>97.6</b>	<b>98.1</b>	<b>93.0</b>	<b>68.5</b>
<i>Simple Texture</i>	Pretrained	96.5	97.2	89.5	54.2
	Random	<b>97.5</b>	<b>97.8</b>	<b>91.9</b>	<b>62.2</b>
<i>Simple Shape</i>	Pretrained	96.9	97.3	92.0	58.8
	Random	<b>97.0</b>	<b>97.6</b>	<b>92.8</b>	<b>62.7</b>

TABLE XI

ABLATION STUDY ON THE LOSSES IN THE T-S FRAMEWORK

Loss	I-AUC	P-AUC	P-PRO	P-mAP
$L_{Cos}$	97.3	97.5	91.1	63.4
$L_{SSIM}$	97.4	98.0	91.5	67.0
DAF	<b>97.6</b>	<b>98.1</b>	<b>93.0</b>	<b>68.5</b>

the *DRA* strategy. It can also be observed that the combination of the discrepancy map  $\bar{M}$  and the segmentation probability map  $M_S$  brings significant increases in both detection and localization performance under the simple strategies.

3) *Influence of the discrepancy map on segmentation*: To investigate the influence of the discrepancy map on the segmentation task, we remove it from the input of the segmentation head *Seg*. As shown in Table VII, the performance suffers from degradation when the discrepancy map is not incorporated. For instance, under the *Simple Texture* strategy, the detection performance drops from 97.5% to 96.5%, while the localization P-mAP reduces by 2.1%. We hypothesize that our improvements lie in the discrepancy map, which provides an effective cue for the segmentation task. This alleviates the segmentation head from the limitation of the synthetic appearance, thereby enhancing its capacity to identify previously unseen anomalies.

4) *Influence of the hyperparameter in the synthesis strategy*: Following DRÆM [10], we evaluate how the hyperparameter  $\beta$  impacts the performance of our method. The result is reported in Table VIII, where DTD and Perlin noise are utilized to simulate the appearance and shape of anomalies in the *DRA* strategy.  $\beta$  controls the opacity in blending. Compared with DRÆM, our method is more robust on  $\beta$ . We keep the same P-mAP as DRÆM, but show advantages in I-AUC (97.6% vs. 97.4%) and P-AUC (97.8% vs. 95.0%) when training via strategy *DRA* without  $\beta$ .

5) *Model Ensemble vs. End-to-End*: Table IX reports the performance of the model ensemble (*Ensem.*). Here we train the teacher-student (*T-S*), and a segmentation model with auxiliary heads alone, then add the discrepancy and the segmentation probability map for evaluation. Compared with *Ensem.*, we achieve improvements in both anomaly detection and localization. The result suggests that with the guidance of the discrepancy map, our end-to-end framework could distinguish normal and abnormal regions more accurately.

6) *Influence of initialization for the student*: Table X studies the impact of parameter initialization on the student *S*. As shown, the student initialized randomly performs better, indicating that pre-trained parameters hinder *S* from yielding

disparate representations from the teacher on anomalies. We assume that the discrepancy map derived from  $T$ - $S$  provides only limited cue for the segmentation head  $Seg$  in this case, resulting in reduced robustness of  $Seg$ .

7) *Influence of the distillation loss*: Table XI shows the effectiveness of the structural similarity ( $SSIM$ ) loss. Since the cosine similarity loss solely distills knowledge at each position in isolation, we introduce  $SSIM$  loss to consider neighboring vectors. The study reveals the effectiveness of the  $SSIM$  loss. DAF achieves the best results when the cosine similarity and  $SSIM$  are adopted simultaneously.

8) *Limitations*: While synthesized anomalies boost performance, the training time increases due to the computationally expensive synthesizing steps. Although adopting simple and cheap ones accelerates the process, it is still slower than unsupervised methods that are solely trained on normal images.

9) *Future work*: In future work, investigating lightweight architectures and developing real-time approaches will be important for industrial applications. Moreover, integrating the rich knowledge derived from large models such as CLIP [35] into our method will be helpful for further enhancing the generalization capacity.

## V. CONCLUSION

In this paper, we have observed that the existing self-supervised methods are susceptible to the quality of synthetic data. To improve the robustness of the prior arts, we have proposed a simple yet effective framework, named Discrepancy Aware Framework (DAF). DAF introduces the teacher-student model to yield the discrepancy map and serves it as an additional cue that reveals the location of anomalies to the segmentation decoder, alleviating the decoder's reliance on the appearance of synthetic data. Meanwhile, the complement of the discrepancy map for segmentation contributes significantly to robustness as well. Extensive experiments have shown that DAF surpasses previous self-supervised methods significantly when faced with simple synthetic strategies in anomaly detection and localization, demonstrating its excellent robustness.

## REFERENCES

- [1] X. Ni, Z. Ma, J. Liu, B. Shi, and H. Liu, "Attention network for rail surface defect detection via consistency of intersection-over-union(iou)-guided center-point estimation," *IEEE Trans. Ind. Informat.*, vol. 18, pp. 1694–1705, 2021. 1
- [2] D. Carrera, F. Manganini, G. Boracchi, and E. Lanzarone, "Defect detection in sem images of nanofibrous materials," *IEEE Trans. Ind. Informat.*, vol. 13, no. 2, pp. 551–561, 2017. 1
- [3] D. Gong, L. Liu, V. Le, B. Saha, M. R. Mansour, S. Venkatesh, and A. van den Hengel, "Memorizing normality to detect anomaly: Memory-augmented deep autoencoder for unsupervised anomaly detection," in *Proc. ICCV*, 2019, pp. 1705–1714. 1
- [4] Y. Cai, H. Chen, X. Yang, Y. Zhou, and K. Cheng, "Dual-distribution discrepancy for anomaly detection in chest x-rays," in *Proc. MICCAI*, 2022, pp. 584–593. 1
- [5] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger, "Uninformed students: Student-teacher anomaly detection with discriminative latent embeddings," in *Proc. CVPR*, 2020, pp. 4183–4192. 1, 3, 6
- [6] G. Wang, S. Han, E. Ding, and D. Huang, "Student-teacher feature pyramid matching for anomaly detection," in *Proc. BMVC*, 2021, pp. 1–14. 1, 3, 4, 6, 7
- [7] H. Deng and X. Li, "Anomaly detection via reverse distillation from one-class embedding," in *Proc. CVPR*, 2022, pp. 9737–9746. 1, 3, 4, 6, 7
- [8] Z. You, L. Cui, Y. Shen, K. Yang, X. Lu, Y. Zheng, and X. Le, "A unified model for multi-class anomaly detection," in *Proc. NeurIPS*, 2022, pp. 4571–4584. 1
- [9] C.-L. Li, K. Sohn, J. Yoon, and T. Pfister, "Cutpaste: Self-supervised learning for anomaly detection and localization," in *Proc. CVPR*, 2021, pp. 9664–9674. 1, 3, 6
- [10] V. Zavrtnik, M. Kristan, and D. Skočaj, "Draem - a discriminatively trained reconstruction embedding for surface anomaly detection," in *Proc. ICCV*, October 2021, pp. 8330–8339. 1, 2, 3, 5, 6, 7, 8
- [11] H. M. Schlüter, J. Tan, B. Hou, and B. Kainz, "Natural synthetic anomalies for self-supervised anomaly detection and localization," in *Proc. ECCV*, 2022. 1, 2, 3, 5, 6, 7, 8
- [12] P. Pérez, M. Gangnet, and A. Blake, "Poisson image editing," in *ACM SIGGRAPH 2003 Papers*, 2003, pp. 313–318. 2
- [13] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger, "Mvtec ad—a comprehensive real-world dataset for unsupervised anomaly detection," in *Proc. CVPR*, 2019, pp. 9592–9600. 2, 5, 6, 8
- [14] M. Wieler and T. Hahn, "Weakly supervised learning for industrial optical inspection," in *DAGM symposium in*, 2007. 2, 5, 7
- [15] Željko Hocenski, S. Vasilic, and V. Hocenski, "Improved canny edge detector in ceramic tiles defect detection," *IECON 2006 - 32nd Annual Conference on IEEE Industrial Electronics*, pp. 3328–3331, 2006. 2
- [16] W.-C. Li and D. ming Tsai, "Defect inspection in low-contrast led images using hough transform-based nonstationary line detection," *IEEE Trans. Ind. Informat.*, vol. 7, pp. 136–147, 2011. 2
- [17] W. Liu, R. Li, M. Zheng, S. Karanam, Z. Wu, B. Bhanu, R. J. Radke, and O. Camps, "Towards visually explaining variational autoencoders," in *Proc. CVPR*, 2020, pp. 8642–8651. 2
- [18] S. Venkataramanan, K.-C. Peng, R. V. Singh, and A. Mahalanobis, "Attention guided anomaly localization in images," in *Proc. ECCV*. Springer, 2020, pp. 485–503. 2
- [19] C. Baur, B. Wiestler, S. Albarqouni, and N. Navab, "Deep autoencoding models for unsupervised anomaly segmentation in brain mr images," in *International MICCAI Brainlesion Workshop*. Springer, 2018, pp. 161–169. 2
- [20] H. Yang, Q. Zhou, K. Song, and Z. Yin, "An anomaly feature-editing-based adversarial network for texture defect visual inspection," *IEEE Trans. Ind. Informat.*, vol. 17, pp. 2220–2230, 2021. 2
- [21] T. Schlegl, P. Seeböck, S. M. Waldstein, G. Langs, and U. Schmidt-Erfurth, "f-anogan: Fast unsupervised anomaly detection with generative adversarial networks," *Medical Image Analysis*, vol. 54, pp. 30–44, 2019. 2
- [22] K. He, X. Chen, S. Xie, Y. Li, P. Doll'ar, and R. B. Girshick, "Masked autoencoders are scalable vision learners," in *Proc. CVPR*, 2021, pp. 15 979–15 988. 2
- [23] K. He, H. Fan, Y. Wu, S. Xie, and R. B. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proc. CVPR*, 2019, pp. 9726–9735. 2
- [24] V. Zavrtnik, M. Kristan, and D. Skočaj, "Reconstruction by inpainting for visual anomaly detection," *Pattern Recognition*, vol. 112, p. 107706, 2021. 3
- [25] N. Ristea, N. Madan, R. T. Ionescu, K. Nasrollahi, F. S. Khan, T. B. Moeslund, and M. Shah, "Self-supervised predictive convolutional attentive block for anomaly detection," in *Proc. CVPR*, 2022, pp. 13 566–13 576. 3
- [26] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Proc. MICCAI*. Springer, 2015, pp. 234–241. 3
- [27] Y. Zou, J. Jeong, L. Pemula, D. Zhang, and O. Dabeer, "Spot-the-difference self-supervised pre-training for anomaly detection and segmentation," in *Proc. ECCV*. Springer, 2022, pp. 392–408. 3
- [28] M. Salehi, N. Sadjadi, S. Baselizadeh, M. H. Rohban, and H. R. Rabiee, "Multiresolution knowledge distillation for anomaly detection," in *Proc. CVPR*, 2021, pp. 14 902–14 912. 3, 4, 5, 6
- [29] G. Tong, Q. Li, and Y. Song, "Two-stage reverse knowledge distillation incorporated and self-supervised masking strategy for industrial anomaly detection," *Knowledge-Based Systems*, vol. 273, p. 110611, 2023. 3
- [30] X. Jiang, J. Liu, J. Wang, Q. Nie, K. Wu, Y. Liu, C. Wang, and F. Zheng, "Softpatch: Unsupervised anomaly detection with noisy data," in *Proc. NeurIPS*, 2022, pp. 15 433–15 445. 3
- [31] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proc. CVPR*, 2009, pp. 248–255. 3
- [32] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. CVPR*, 2016, pp. 770–778. 4

- [33] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi, "Describing textures in the wild," in *Proc. CVPR*, 2014, pp. 3606–3613. [5](#), [7](#)
- [34] J. Božič, D. Tabernik, and D. Skočaj, "End-to-end training of a two-stage neural network for defect detection," in *Proc. ICPR*. IEEE, 2021, pp. 5619–5626. [6](#)
- [35] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *Proc. ICML*, 2021, pp. 8748–8763. [9](#)