

# Revisiting Deep Feature Reconstruction for Logical and Structural Industrial Anomaly Detection

Anonymous authors

Paper under double-blind review

## Abstract

Industrial anomaly detection is crucial for quality control and predictive maintenance but is challenging due to limited training data, varied anomaly types, and changing external factors affecting object appearances. Existing methods detect structural anomalies, such as dents and scratches, by relying on multi-scale features of image patches extracted from a deep pre-trained network. Nonetheless, extensive memory or computing requirement hinders their adoption in practice. Furthermore, detecting logical anomalies, such as images with missing or surplus elements, necessitates understanding spatial relationships beyond traditional patch-based methods. Our work focuses on Deep Feature Reconstruction (DFR), which offers a memory- and compute-efficient way of detecting structural anomalies. Moreover, we extend DFR to develop a unified framework for detecting structural and logical anomalies, called ULSAD. Specifically, we improve the training objective of DFR to enhance the capability to detect structural anomalies and introduce an attention-based loss using a global autoencoder-like network for detecting logical anomalies. Empirical results on five benchmark datasets demonstrate the effectiveness of ULSAD in the detection and localization of both structural and logical anomalies compared to eight state-of-the-art approaches. Moreover, an in-depth ablation study showcases the importance of each component in enhancing overall performance. Our code can be accessed here: <https://anonymous.4open.science/r/ULSAD-2024>.

## 1 Introduction

Anomaly detection (AD) is a widely studied problem in machine learning that is used to identify rare events or unusual patterns (Salehi et al., 2022). It enables the detection of abnormalities, potential threats, or critical system failures across diverse applications such as predictive maintenance (PdM) (Tang et al., 2020; Choi et al., 2022), fraud detection (Ahmed et al., 2016; Hilal et al., 2022), and medicine (Tibshirani & Hastie, 2007; Fernando et al., 2021). Despite its importance and widespread applicability, it remains a challenging task as the anomalous samples are not known a priori (Ruff et al., 2021). Typically, AD is therefore addressed as an unsupervised representation learning problem (Pang et al., 2020) where the training data contains predominantly normal samples. Therefore, the aim is to learn the normal behaviour using the samples in the training set and identify anomalies as deviations from this normal behaviour. This setting is also known as one-class classification (Ruff et al., 2018).

Our study concentrates on Industrial Anomaly Detection (IAD) (Bergmann et al., 2019), specifically targeting the detection and localization of anomalies in images from industrial manufacturing processes. Over the years, it has garnered attention in both industry and academia as AD can be used for various tasks like quality control or predictive maintenance, which are of primal interest to industries. Despite the necessity, addressing IAD is difficult given the following challenges: (i) evolving processes resulting in different manifestations of anomaly (Gao et al., 2023), and (ii) varying object appearances due to external factors such as background, lighting conditions and orientation (Jezek et al., 2021). Furthermore, the anomalies in IAD can be broadly categorized into: (i) **structural anomalies** where subtle localized structural defects can be observed (Bergmann et al., 2019), and (ii) **logical anomalies** where violations of logical constraints result in anomalies (Bergmann et al., 2022). Examples of both structural and logical anomalies can be seen in Figure 1.

Methods proposed for the detection of **structural anomalies** leverage multi-scale features of image patches obtained using deep convolutional neural networks (Salehi et al., 2021). PatchCore (Roth et al., 2022) and

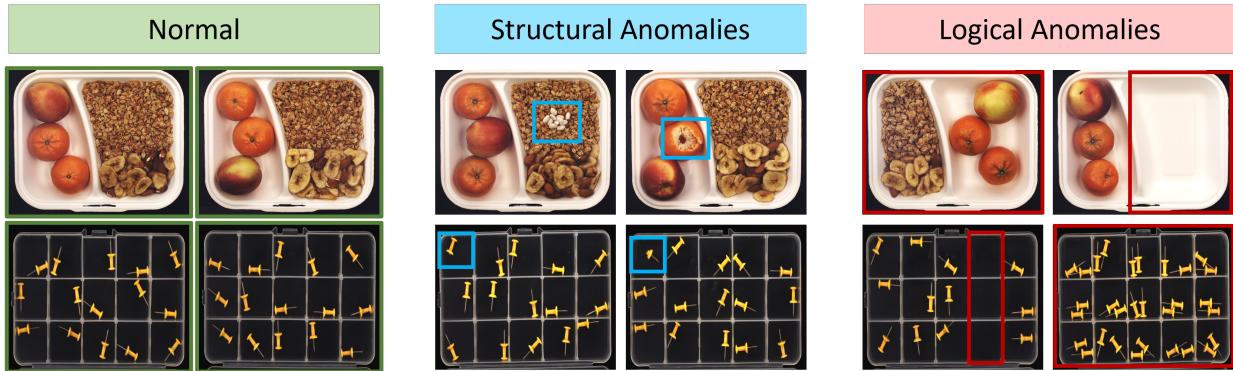


Figure 1: Types of anomalies. (Left) First, two normal samples are shown from the categories “breakfast box” and “pushpin” of the MVTecLOCO dataset. Then, we show examples of structural (Middle) and logical anomalies (Right) with the anomalies highlighted in blue and red, respectively.

CFA(Lee et al., 2022) achieved state-of-the-art (SOTA) performance by storing the extracted features in a memory bank and comparing the features of the test image with their closest neighbour from the memory bank. However, such approaches require considerable storage to accumulate the extracted features, which can be challenging for large-scale datasets. The first alternative is knowledge distillation-based approaches (Bergmann et al., 2020; 2022), where a student network is trained to mimic the teacher for normal samples. During inference, anomalies are identified based on the discrepancy between the student and teacher output. A key requirement of these distillation-based approaches is that the student network must be less expressive than the teacher to prevent it from mimicking the teacher on anomalous samples. Thus, regularization methods, such as penalty based on an external dataset or hard-mining loss (Batzner et al., 2024) are applied, which slows down the training and increases the requirement of computing resources. Moreover, excessive regularization also prevents learning representations for normal images. The second alternative is to model the features of normal images using a multivariate Gaussian distribution (Defard et al., 2021) or learn to reconstruct the features using a deep feature reconstruction (DFR) network (Yang et al., 2020).

Besides structural anomalies, **logical anomalies** occur when elements in the images are missing, misplaced, in surplus or violate geometrical constraints (Bergmann et al., 2022). Methods relying on multi-scale features of image patches would fail as they would still be considered normal. It is the combination of objects in the image that makes the image anomalous. Thus, to detect such logical anomalies, it is necessary to look beyond image patches and develop a global understanding of the spatial relationships within normal images. Distillation-based methods, which are predominantly used for the detection of logical anomalies, rely on an additional network to learn the spatial relationships between items in the normal image (Batzner et al., 2024).

In this paper, we focus on DFR, the benefits of which are four-fold. First, it does not need large memory for storing the features, unlike PatchCore (Roth et al., 2022). Second, unlike PaDiM (Defard et al., 2021), it does not make any assumption about the distribution of features. Third, learning to reconstruct features in the latent space of a pre-trained network is less impacted by the curse of dimensionality than learning to reconstruct high-dimensional images. Fourth, deep networks trained to reconstruct normal images using the per-pixel distance suffer from the loss of sharp edges of the objects or textures in the background. As a consequence, AD performance deteriorates due to an increase in false positives, i.e., the number of normal samples falsely labelled anomalies. On the contrary, using the distance between the multi-scale features and their corresponding reconstructions during training is less likely to result in such errors (Assran et al., 2023).

We revisit DFR to develop a unified framework for the detection of structural and logical anomalies. First, to improve the performance of DFR for the detection of structural anomalies, we modify the training objective by considering a combination of  $\ell_2$  and cosine distances between each feature and the corresponding reconstruction. The incorporation of the cosine distance addresses the curse of dimensionality, where high-dimensional features become uniformly distant from each other in Euclidean space (Aggarwal et al., 2001). Second, to allow for the detection of logical anomalies, we introduce an attention-based loss using a global autoencoder-like network. We empirically demonstrate that with our proposed improvements, not only do

the detection and localization capabilities of DFR improve for structural anomalies, but also it delivers competitive results on the detection of logical anomalies. Our contributions can be summarized as:

- We propose **Unified framework for Logical and Structural Anomaly Detection** referred as ULSAD, a framework for identifying and localizing both structural and logical anomalies building on DFR.
- We improve DFR for structural AD by considering the angular difference between the extracted and reconstructed feature vectors in addition to the difference in values.
- For learning the logical constraints, we propose a novel attention-based loss for learning the features of an autoencoder which improves the detection and localization of logical anomalies during inference.
- We demonstrate the effectiveness of ULSAD by comparing it with 8 SOTA methods across 5 widely adopted IAD benchmark datasets.
- Through extensive ablation study, we show the effect of each component of ULSAD on the overall performance of the end-to-end architecture.

## 2 Related Work

Several methods have been proposed over the years for addressing Industrial AD (Bergmann et al., 2019; 2022; Jezek et al., 2021). They can be broadly categorized into feature-embedding based and reconstruction-based methods. We briefly highlight some relevant works in each of the category. For an extended discussion on the prior works we refer the readers to the survey by Liu et al. (2024).

**Feature Embedding-based** There are mainly three different types of IAD methods which utilize feature embeddings from a pre-trained deep neural network: memory bank (Defard et al., 2021; Roth et al., 2022; Lee et al., 2022), student-teacher (Zhang et al., 2023; Batzner et al., 2024), and density-based (Gudovskiy et al., 2021; Yu et al., 2021). The main idea of memory bank methods is to extract features of nominal images and store them in a memory bank during the training phase. During the testing phase, the feature of a test image is used as a query to match the stored nominal features. There are two main constraints in these methods: *how to learn useful features* and *how to reduce the size of the memory bank*. PatchCore (Roth et al., 2022) introduces a coreset selection algorithm to reduce the memory bank size. CFA (Lee et al., 2022) proposes contrastive supervision based on a coupled hypersphere to learn target-oriented features and a compression scheme to reduce memory size. The performance of the memory bank methods heavily depends on the completeness of the memory bank, which requires a large number of nominal images. Moreover, the memory size is often related to the number of training images, which makes these methods not preferable for large datasets or very high-dimensional images. In the student-teacher approach, the student network learns to extract features of the nominal samples, similar to the teacher model. For anomalous images, the features extracted by the student network should be different from the teacher network. Batzner et al. (2024) propose to use an autoencoder model in addition to the student network to identify logical anomalies. For leveraging the multiscale feature from the teacher network to detect anomalies at various scales, Deng & Li (2022) propose **Reverse Distillation**. Zhang et al. (2023) extended it by proposing to utilize two student networks to deal with structural and logical anomalies. Yang et al. (2020) propose to learn a deep neural network for learning to reconstruct the features of the normal images extracted using the pre-trained backbone. On the other hand, density-based methods detect anomalies based on the likelihood of an extracted feature of the test sample, given the estimated distribution of the features of the normal samples. PaDiM (Defard et al., 2021) uses a multivariate Gaussian to learn the probabilistic representation of the nominal class while FastFlow (Yu et al., 2021) and CFLOW (Gudovskiy et al., 2021) utilize normalizing flows.

**Reconstruction-based methods.** Reconstruction-based methods assume that encoder-decoder models trained on normal samples will exhibit poor performance for anomalous samples. However, relying solely on the reconstruction objective can result in the model collapsing to an identity mapping. To address this, structural assumptions are made regarding the data generation process. One such assumption is the Manifold Assumption, which posits that the observed data resides in a lower-dimensional manifold within the data space. Methods leveraging this assumption impose a bottleneck by restricting the encoded space to a lower dimensionality than the actual data space. Common deep reconstruction models used include

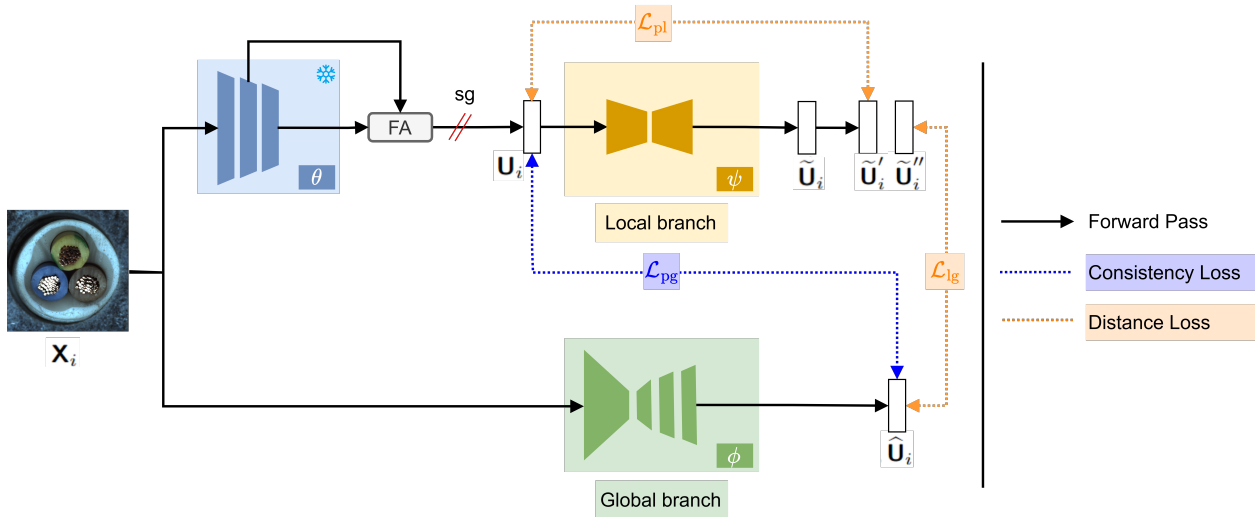


Figure 2: Overview of the end-to-end architecture of ULSAD

AE or VAE-based approaches. Advanced strategies encompass techniques like reconstruction by memorised normality (Gong et al., 2019), model architecture adaptation (Lai et al., 2019) and partial/conditional reconstruction (Yan et al., 2021; Nguyen et al., 2019). Generative models like GANs are also widely employed for anomaly detection, as the discriminator inherently calculates reconstruction loss for samples (Zenati et al., 2018). Variants of GANs, such as denoising GANs (Sabokrou et al., 2018) and class-conditional GANs (Perera et al., 2019), improve anomaly detection performance by increasing the challenge of reconstruction. Some methods utilize the reconstructed data from GANs in downstream tasks to enhance the amplification of reconstruction errors for anomaly detection (Zhou et al., 2020). Lastly, DRÆM (Zavrtanik et al., 2021) trains an additional discriminative network alongside a reconstruction network to improve the AD performance.

In this paper, we focus on feature embedding-based methods motivated by their effectiveness in the current SOTA methods. Specifically, we build on DFR (Yang et al., 2020), which has several benefits. First, it is memory-efficient as it does not rely on a memory bank of extracted features, unlike (Roth et al., 2022). Second, unlike (Defard et al., 2021), it does not make any assumption about the distribution of the extracted features. Third, it is computationally efficient and less impacted by the curse of dimensionality as it operates in the lower-dimensional latent space of a deep neural network. Last, by avoiding the use of per-pixel distance in its reconstruction objective, it is less prone to false positives (Assran et al., 2023).

### 3 The ULSAD Framework for Anomaly Detection

We propose ULSAD, a framework for simultaneously detecting and localising anomalies in images as shown in Figure 2. Firstly, we utilize a feature extractor network for extracting low-dimensional features from high-dimensional images, which we discuss in Section 3.1. Then, for the detection of structural and logical anomalies, we rely on a dual-branch architecture. The local branch detects structural anomalies with the help of a feature reconstruction network applied to the features corresponding to patches in the input image. We elaborate on this in Section 3.2. Conversely, the global branch, as discussed in Section 3.3, detects logical anomalies using an autoencoder-like network, which takes as input the entire image. Lastly, we provide an overview of the training algorithm of ULSAD in Section 3.4 followed by a discussion on the inference process in Section 3.5.

We consider a dataset  $\mathcal{D} = \{(\mathbf{X}_i, y_i)\}_{i=1}^n$  with  $n$  samples where  $\mathbf{X}_i \in \mathcal{X}$  is an input and  $y_i \in \mathcal{Y} := \{0, 1\}$  is the corresponding label. It is important to note that for this work, we specifically focus on image-based AD where each input  $\mathbf{X}_i$  is an image. We refer to the normal class with the label 0 and the anomalous class with the label 1. The samples belonging to the anomalous class can contain either logical or structural anomalies or a combination of both. We denote the disjoint train and test partitions of  $\mathcal{D}$  as  $\mathcal{D}_{\text{train}}$  and  $\mathcal{D}_{\text{test}}$ . The training set contains only normal samples, i.e.,  $y_i = 0$ . Therefore, for the sake of simplicity, we refer to it as  $\mathcal{D}_N = \{\mathbf{X}_i | (\mathbf{X}_i, y_i) \in \mathcal{D}_{\text{train}}\}$ . The test set  $\mathcal{D}_{\text{test}}$  includes both normal and anomalous samples.

### 3.1 Feature Extractor

High-dimensional images pose a significant challenge for AD (Reiss et al., 2022). Recent studies have shown that deep convolutional neural networks (CNNs) trained on ImageNet (Russakovsky et al., 2015) capture discriminative features for several downstream tasks. Typically, AD methods (Salehi et al., 2021; Defard et al., 2021; Yoon et al., 2023) leverage such pre-trained networks to extract features corresponding to partially overlapping regions or patches in the images. **Learning to detect anomalies using the lower-dimensional features is beneficial as it results in reduced computational complexity.** A key factor determining the efficiency of such methods is the size of the image patches being used, as anomalies can occur at any scale. To overcome this challenge, features are extracted from multiple layers of the CNNs and fused together (Salehi et al., 2021; Roth et al., 2022; Yang et al., 2020). As features from each layer have a different receptive field, each element of the extracted feature would correspond to a patch of a different size in the image. Thus, combining features from multiple layers results in multi-scale features of image patches, referred to as **patch feature**.

Similar to DFR, we extract low-dimensional feature maps by combining features from multiple layers of a pre-trained network  $\theta$ . **In this paper, we consider ResNet-like architectures for  $\theta$ .** With the increasing number of layers, the computation becomes increasingly expensive as the resulting tensor becomes high-dimensional. In order to overcome this, we consider two intermediate or mid-level features. Our choice is guided by the understanding that the initial layers of such deep networks capture generic image features, while the latter layers are often biased towards the pre-training classification task (Roth et al., 2022). We denote the features extracted at a layer  $j$  for an image  $\mathbf{X}_i$  as  $\theta_j(\mathbf{X}_i)$ . Following this convention, we express the feature map  $\mathbf{U}_i \in \mathcal{U} = \mathbb{R}^{c^* \times h^* \times w^*}$  produced by the *Feature Aggregator* (FA) as a concatenation of  $\theta_j(\mathbf{X}_i)$  and  $\theta_{j+1}(\mathbf{X}_i)$  obtained from layers  $j$  and  $j+1$  of  $\theta$ . Furthermore, to facilitate the concatenation of features extracted from multiple layers of the extractor  $\theta$ , the features at the lower resolution layer  $j+1$  are linearly rescaled by FA to match the dimension of the features at layer  $j$ . Furthermore, we define an invertible transformation  $f: \mathbb{R}^{c^* \times h^* \times w^*} \rightarrow \mathbb{R}^{c^* \times (h^* \times w^*)}$ . Now, using  $f$ , we express  $\mathbf{Z}_i = f(\mathbf{U}_i)$ . Moreover,  $\mathbf{Z}_i[:, k] = \mathbf{U}_i[:, h, w]$  with  $k = (h-1) \times w^* + w$  where  $h \in \{1, 2, \dots, h^*\}$ ,  $w \in \{1, 2, \dots, w^*\}$ , and  $k \in \{1, 2, \dots, h^* \times w^*\}$ . The function  $f$  can be computed in practice by reshaping the tensor to obtain a 2D matrix.

### 3.2 Detecting Structural Anomalies

Given a transformed feature map  $\mathbf{Z}_i$  as described in the previous section, we elaborate on the local branch of ULSAD for the detection of structural anomalies. Specifically, our goal is to train a neural network to reconstruct the patch features given the dataset  $\mathcal{D}_N$  comprising solely of normal images. Therefore, we can identify the structural anomalies when the network fails to reconstruct a patch feature during inference.

We refer to each vector in the feature map  $\mathbf{Z}_i$  as  $\mathbf{z}_{ik} = \mathbf{Z}_i[:, k] \in \mathbb{R}^{c^* \times 1}$ . Now, recall that  $\theta$  is a CNN with multiple convolutional and pooling layers. Therefore, each  $c^*$ -dimensional feature vector  $\mathbf{z}_{ik}$  of the output feature map  $\mathbf{Z}_i$  has a receptive field greater than one. In other words, each feature vector  $\mathbf{z}_{ik}$  corresponds to a patch on the original image  $\mathbf{X}_i$  and hence can be considered as a patch feature. The size of the patch in the original image is determined by the receptive field of  $\theta$ .

**Feature Reconstruction Network (FRN).** ULSAD utilizes a convolutional encoder-decoder architecture with a lower-dimensional bottleneck for learning to reconstruct the feature map  $\mathbf{U}_i$  given the normal images in training dataset  $\mathcal{D}_N$ , as illustrated in Figure 3. First, the encoder network  $\psi_e$  compresses the feature  $\mathbf{U}_i$  to a lower dimensional space, which induces the information bottleneck. It acts as an implicit regularizer, preventing generalization to features corresponding to anomalous images. The encoded representation is then mapped back to the latent space using a decoder network  $\psi_d$ . Therefore, we can express the FRN as  $\psi = \psi_e \circ \psi_d$ . The output of FRN is  $\tilde{\mathbf{U}}_i = \psi(\mathbf{U}_i) \in \mathbb{R}^{2c^* \times h^* \times w^*}$ . Note that the number of channels is doubled to simultaneously generate two feature maps  $\tilde{\mathbf{U}}_i'$  and  $\tilde{\mathbf{U}}_i''$ , both having dimension  $c^* \times h^* \times w^*$ . For the local objective, we utilize  $\tilde{\mathbf{U}}_i'$ , while  $\tilde{\mathbf{U}}_i''$  is used in the global objective. **Although the feature maps**

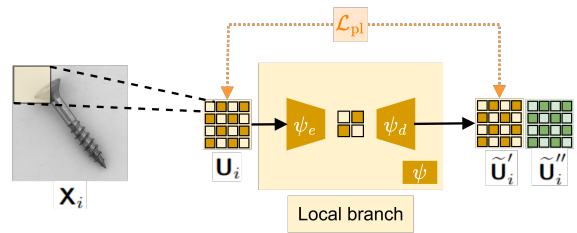


Figure 3: Feature Reconstruction Network

$\mathbf{U}_i$  have significantly lower dimensionality compared to the input images  $\mathbf{X}_i$ , they can still be considered high-dimensional tensors. In high-dimensional spaces, the  $\ell_2$  distance is not effective at distinguishing between the nearest and furthest points (Aggarwal et al., 2001), making it an inadequate measure for computing the difference between feature maps during training. Therefore, similar to Salehi et al. (2021), we propose combining  $\ell_2$  and cosine distances to account for differences in both the value and direction of the vectors in the feature maps as follows:

$$\mathcal{L}_{\text{pl}}(\tilde{\mathbf{Z}}'_i, \mathbf{Z}_i) = \frac{1}{h^* \times w^*} \sum_{k=1}^{h^* w^*} l_v(\tilde{\mathbf{z}}'_{ik}, \mathbf{z}_{ik}) + \lambda_l l_d(\tilde{\mathbf{z}}'_{ik}, \mathbf{z}_{ik}), \quad (1)$$

where  $\tilde{\mathbf{Z}}'_i = f(\tilde{\mathbf{U}}'_i)$ ,  $\mathbf{Z}_i = f(\mathbf{U}_i)$  and  $\lambda_l \geq 0$  is a hyperparameter. Furthermore,  $l_v(\tilde{\mathbf{z}}'_{ik}, \mathbf{z}_{ik})$  measures the difference in values between the vectors  $\mathbf{z}_{ik}$  and  $\tilde{\mathbf{z}}'_{ik}$ , while  $l_d(\tilde{\mathbf{z}}'_{ik}, \mathbf{z}_{ik})$  measures the angular distance between them. We define  $l_v$  and  $l_d$  as

$$l_v(\tilde{\mathbf{z}}'_{ik}, \mathbf{z}_{ik}) = \frac{1}{c^*} \|\tilde{\mathbf{z}}'_{ik} - \mathbf{z}_{ik}\|_2^2, \quad \text{and} \quad l_d(\tilde{\mathbf{z}}'_{ik}, \mathbf{z}_{ik}) = 1 - \frac{(\tilde{\mathbf{z}}'_{ik})^T \mathbf{z}_{ik}}{\|\tilde{\mathbf{z}}'_{ik}\|_2 \|\mathbf{z}_{ik}\|_2}. \quad (2)$$

### 3.3 Detecting Logical Anomalies

Although the feature reconstruction task discussed in Section 3.2 allows us to detect structural anomalies, it is not suited for identifying logical anomalies that violate the logical constraints of normal images. Recall that such violations appear in the form of misplaced, misaligned, or surplus objects found in normal images. If we consider the example of misaligned objects, the previously discussed approach will fail as it focuses on the individual image patches, which would be normal. It is the overall spatial arrangement of objects in the image which is anomalous. Thus, to identify such anomalies, our goal is to learn the spatial relationships among the objects present in the normal images of the training dataset  $\mathcal{D}_N$ . We achieve this with the global branch of ULSAD, shown in Figure 4, which leverages the entire image and not just its individual patches.

In order to achieve our goal, we start with an observation of the features extracted using the pre-trained network  $\theta$ . Pre-trained CNNs tend to have similar activation patterns for semantically similar objects (Tung & Mori, 2019; Zagoruyko & Komodakis, 2017). In Figure 5, we visualize four attention maps computed from the features of a pre-trained Wide-Resnet50-2 network. It can be seen that in the first map, all the items for the semantic class “fruits” receive a high attention score. The remaining attention maps focus on individual semantic concepts like “oranges”, “cereal” and “plate”, respectively.

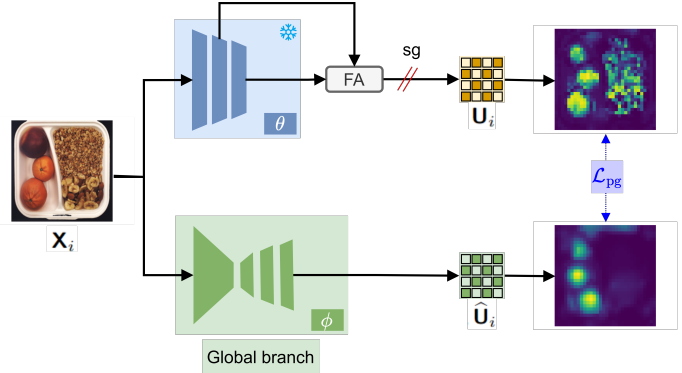


Figure 4: Global Branch of ULSAD

Based on this observation and inspired by the attention-transfer concept for knowledge distillation (Zagoruyko & Komodakis, 2017; Tung & Mori, 2019), we propose to learn the spatial relationships (Dosovitskiy et al., 2021) among the feature vectors in  $\mathbf{U}_i$  obtained from normal images. Recall that each feature vectors correspond to a patch in the original image. Therefore, learning the spatial relationships among the feature vectors would allow us to learn the same among the patches in the original image. **This forces ULSAD to learn the relative positions of objects in the normal images, thereby enabling it to capture the logical constraints.** Starting from  $\mathbf{Z}_i = f(\mathbf{U}_i)$ , we first compute the self-attention weight matrix  $\mathbf{W}_i \in \mathbb{R}^{(h^* \times w^*) \times (h^* \times w^*)}$ , such that

$$\mathbf{W}_i[p, q] = \frac{\exp(\mathbf{z}_{ip}^T \mathbf{z}_{iq} / \sqrt{c^*})}{\sum_{r=1}^{h^* \times w^*} \exp(\mathbf{z}_{ir}^T \mathbf{z}_{iq} / \sqrt{c^*})}. \quad (3)$$

Then, the attention map  $\mathbf{A}_i \in \mathbb{R}^{c^* \times (h^* \times w^*)}$  is computed as  $\mathbf{A}_i = \mathbf{Z}_i \mathbf{W}_i$ . For learning the spatial relations using  $\mathbf{A}_i$  as our target, we use a convolutional autoencoder-like network  $\phi = \phi_e \circ \phi_d$  where  $\phi_e$  is the encoder

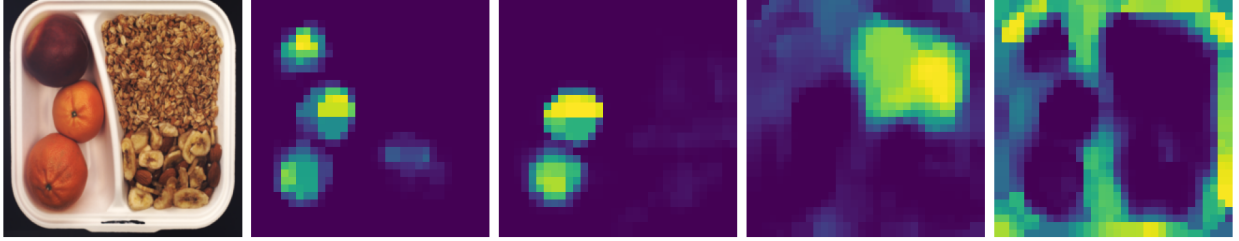


Figure 5: (First) Example image belonging to the category ‘‘breakfast box’’ in the MVTeCLOCO dataset. (Rest) Visualization of attention maps computed using the intermediate features from a pre-trained model.

and  $\phi_d$  is the decoder. Similar to a standard autoencoder,  $\phi_e$  compresses the input image  $\mathbf{X}_i$  to a lower dimensional space. However,  $\phi_d$  maps the encoded representation to the feature space  $\mathcal{U}$ , which has a lower dimension than the input space  $\mathcal{X}$ . We denote the output of  $\phi$  as  $\hat{\mathbf{U}}_i = \phi(\mathbf{X}_i)$ .

A direct approach would be to compute the self-attention map for  $\hat{\mathbf{U}}_i$  and minimize its distance from  $\mathbf{A}_i$ . However, it makes the optimization problem computationally challenging as each vector in  $\hat{\mathbf{U}}_i$  is coupled with every other vector by the network weights  $\phi$  (Zhang et al., 2023). To overcome this, we compute the cross-attention map  $\hat{\mathbf{A}}_i \in \mathbb{R}^{c^* \times (h^* \times w^*)}$  between  $\mathbf{U}_i$  and  $\hat{\mathbf{U}}_i$ , by first computing  $\widehat{\mathbf{W}}_i$  using  $\hat{\mathbf{Z}}_i = f(\hat{\mathbf{U}}_i)$  as

$$\widehat{\mathbf{W}}_i[p, q] = \frac{\exp(\mathbf{z}_{ip}^T \hat{\mathbf{z}}_{iq} / \sqrt{c^*})}{\sum_{r=1}^{h^* \times w^*} \exp(\mathbf{z}_{ir}^T \hat{\mathbf{z}}_{iq} / \sqrt{c^*})}. \quad (4)$$

Then, the attention map  $\hat{\mathbf{A}}_i$  can be computed as  $\hat{\mathbf{A}}_i = \mathbf{Z}_i \widehat{\mathbf{W}}_i$ . Given, the self-attention map  $\mathbf{A}_i$  and the cross-attention map  $\hat{\mathbf{A}}_i$ , the global loss is defined as

$$\mathcal{L}_{\text{pg}}(\hat{\mathbf{A}}_i, \mathbf{A}_i) = \frac{1}{h^* \times w^*} \sum_{k=1}^{h^* \times w^*} l_v(\hat{\mathbf{a}}_{ik}, \mathbf{a}_{ik}) + \lambda_g l_d(\hat{\mathbf{a}}_{ik}, \mathbf{a}_{ik}), \quad (5)$$

where  $\mathbf{a}_{ik} = \mathbf{A}_i[:, k]$ ,  $\hat{\mathbf{a}}_{ik} = \hat{\mathbf{A}}_i[:, k]$  and  $\lambda_g \geq 0$  is a hyperparameter. It is important to note that the loss in Equation 5 is minimized when the attention map  $\hat{\mathbf{A}}_i$  is the mean over all the training samples. As a result, the global branch is prone to false positives in the presence of sharp edges or heavily textured surfaces due to the loss of high-frequency details (Dosovitskiy & Brox, 2016; Assran et al., 2023). To address this, we utilize the FRN  $\psi$  in the local branch to learn the output  $\tilde{\mathbf{U}}_i$ . Recall that the output of FRN  $\tilde{\mathbf{U}} \in \mathbb{R}^{2c^* \times h^* \times w^*}$  has  $2c^*$  number of channels to simultaneously generate two feature maps  $\tilde{\mathbf{U}}'_i$  and  $\tilde{\mathbf{U}}''_i$ , both having dimension  $c^* \times h^* \times w^*$ . Out of which,  $\tilde{\mathbf{U}}'_i$  is used for learning the patch features. Here, we define the loss  $\mathcal{L}_{lg}$  to relate the local feature map  $\tilde{\mathbf{U}}''_i$  with the global feature map  $\hat{\mathbf{U}}_i$  as

$$\mathcal{L}_{lg}(\tilde{\mathbf{Z}}''_i, \hat{\mathbf{Z}}_i) = \frac{1}{h^* \times w^*} \sum_{k=1}^{h^* \times w^*} l_v(\tilde{\mathbf{z}}''_{ik}, \hat{\mathbf{z}}_{ik}) + \lambda_g l_d(\tilde{\mathbf{z}}''_{ik}, \hat{\mathbf{z}}_{ik}), \quad (6)$$

where  $\tilde{\mathbf{Z}}''_i = f(\tilde{\mathbf{U}}''_i)$ . Therefore, during inference, a difference between the  $\tilde{\mathbf{U}}''_i$  and  $\hat{\mathbf{U}}_i$  indicates the presence of logical anomalies with fewer false positives. The benefits of such a framework are two-fold: (1) it allows for learning the spatial relationships in the normal images while minimizing false positives, and (2) doubling the channels in the decoder allows sharing the encoder architecture, reducing the computational costs.

### 3.4 Algorithmic Overview

An overview of ULSAD is outlined in Algorithm 1, which can simultaneous detect structural and logical anomalies. Firstly, we pass a normal image  $\mathbf{X}_i$  from the training dataset  $\mathcal{D}_N$  through the feature extractor  $\theta$  to obtain multi-scale feature  $\mathbf{U}_i$ . We normalize the features (line 4, Algorithm 1) with the channel-wise mean

**Algorithm 1:** Unified Logical and Structural AD (ULSAD) // Local branch Global branch

---

**Require:** Training dataset  $\mathcal{D}_N$ , Feature extractor  $\theta$ , Feature reconstruction network  $\phi$   
Global autoencoder  $\psi$ , Number of epochs  $e$ , Learning rate  $\eta$ , Pre-trained feature statistics  $\mu, \sigma$

- 1 **for** (epoch  $\in 1, \dots, e$ ) and ( $\mathbf{X}_i \in \mathcal{D}_N$ ) **do**
- 2     Extract normalized features maps using the pre-trained network:
- 3      $\mathbf{U}_i \leftarrow \theta(\mathbf{X}_i)$
- 4      $\mathbf{U}_i \leftarrow (\mathbf{U}_i - \mu) / \sigma$
- 5      $\mathbf{Z}_i \leftarrow f(\mathbf{U}_i)$
- 6     Reconstruct the features maps using the local branch:
- 7      $\tilde{\mathbf{U}}_i \leftarrow \psi(\mathbf{U}_i)$
- 8      $\tilde{\mathbf{Z}}_i \leftarrow f(\tilde{\mathbf{U}}_i)$
- 9     Compute local loss (Eq. 1):
- 10     $l_l \leftarrow \mathcal{L}_{\text{pl}}(\tilde{\mathbf{Z}}_i, \mathbf{Z}_i)$
- 11    Obtain the output of the global autoencoder:
- 12     $\hat{\mathbf{U}}_i \leftarrow \phi(\mathbf{X}_i)$
- 13     $\hat{\mathbf{Z}}_i \leftarrow f(\hat{\mathbf{U}}_i)$
- 14    Compute global loss (Eq. 5):
- 15     $l_g \leftarrow \mathcal{L}_{\text{pg}}(\hat{\mathbf{Z}}_i, \mathbf{Z}_i)$
- 16    Compute local-global loss (Eq. 6):
- 17     $l_{\text{lg}} \leftarrow \mathcal{L}_{\text{lg}}(\hat{\mathbf{Z}}_i, \tilde{\mathbf{Z}}_i'')$
- 18    Compute overall loss:
- 19     $l \leftarrow l_l + l_g + l_{\text{lg}}$
- 20    Update model parameters:
- 21     $\phi \leftarrow \phi - \eta \nabla_{\phi} l$
- 22     $\psi \leftarrow \psi - \eta \nabla_{\psi} l$
- 23 **end**

**Return:**  $\phi, \psi$

---

$\mu$  and standard deviation  $\sigma$  computed over all the features. We do not include this step in Algorithm 1 as the calculation is trivial. Instead, we consider the values  $\mu$  and  $\sigma$  to be given as input parameters for the sake of simplicity. Secondly, we obtain  $\tilde{\mathbf{U}}_i$  by passing  $\mathbf{U}_i$  through the feature reconstruction network  $\psi$  (line 7, Algorithm 1). Recall that,  $\tilde{\mathbf{U}}_i$  has a dimension  $2c^* \times h^* \times w^*$  which can be decomposed into two feature maps  $\tilde{\mathbf{U}}_i'$  and  $\tilde{\mathbf{U}}_i''$  each with a dimension  $c^* \times h^* \times w^*$ . The feature reconstruction loss  $\mathcal{L}_{\text{pl}}$  is then computed between  $\mathbf{Z}_i$  and  $\tilde{\mathbf{Z}}_i'$ , where  $\mathbf{Z}_i = f(\mathbf{U}_i)$  and  $\tilde{\mathbf{Z}}_i' = f(\tilde{\mathbf{U}}_i')$ . Thirdly, we obtain the features  $\hat{\mathbf{U}}_i$  by passing the input sample  $\mathbf{X}_i$  through the autoencoder  $\phi$ . Then for learning the spatial relationships from the normal images, we compute  $\mathcal{L}_{\text{pg}}$  between the self-attention map of  $\mathbf{Z}_i$  and the cross-attention map between  $\mathbf{Z}_i$  and  $\hat{\mathbf{Z}}_i = f(\hat{\mathbf{U}}_i)$  (line 15, Algorithm 1). In the fourth step, we compute the loss  $\mathcal{L}_{\text{lg}}$  between  $\hat{\mathbf{Z}}_i$  and  $\tilde{\mathbf{Z}}_i'' = f(\tilde{\mathbf{U}}_i'')$ . Finally, the model parameters  $\psi$  and  $\phi$  are updated based on the gradient of the total loss (line 21 – 22, Algorithm 1). The end-to-end pipeline is illustrated in Figure 2.

### 3.5 Anomaly Detection and Segmentation

First, the local anomaly map  $\mathbf{M}_i^l \in \mathbb{R}^{h^* \times w^*}$  is computed between the output of the local branch  $\tilde{\mathbf{U}}_i'$  and the feature  $\mathbf{U}_i$ , such that

$$\mathbf{M}_i^l[h, w] = l_v(\tilde{\mathbf{U}}_i'[:, h, w], \mathbf{U}_i[:, h, w]) + \lambda_l l_d(\tilde{\mathbf{U}}_i'[:, h, w], \mathbf{U}_i[:, h, w]), \quad (7)$$



where  $\tilde{\mathbf{U}}'_i = f^{-1}(\tilde{\mathbf{Z}}'_i)$  and  $\mathbf{U}_i = f^{-1}(\mathbf{Z}_i)$ . Similarly, the global anomaly map  $\mathbf{M}_i^g$  is calculated using the output from the global autoencoder  $\tilde{\mathbf{Z}}_i$  and the local reconstruction branch  $\tilde{\mathbf{Z}}'_i$ :

$$\mathbf{M}_i^g[h, w] = l_v(\tilde{\mathbf{U}}''_i[:, h, w], \hat{\mathbf{U}}_i[:, h, w]) + \lambda_g l_d(\tilde{\mathbf{U}}''_i[:, h, w], \hat{\mathbf{U}}_i[:, h, w]), \quad (8)$$

Recall that  $\lambda_l$  and  $\lambda_g$  are hyperparameters. To obtain the final anomaly map  $m_i$  for the image  $x_i$ , the local and global anomaly maps are averaged. As local and global anomaly maps can have different ranges of anomaly scores, prior to averaging, each map is normalized individually. Besides ensuring a similar range of scores for both maps, it prevents noise from one of the maps, making the anomaly detected from the other map indistinguishable. In this work, we follow a quantile-based normalization. The benefit is that it does not make any assumptions about the distribution of anomaly scores.

First, we generate two sets of anomaly maps:  $\mathbb{M}^l = \{\mathbf{M}_i^l \mid \mathbf{X}_i \in \mathcal{D}_{\text{valid}}\}$  and  $\mathbb{M}^g = \{\mathbf{M}_i^g \mid \mathbf{X}_i \in \mathcal{D}_{\text{valid}}\}$  utilizing images from the validation dataset  $\mathcal{D}_{\text{valid}}$ . Subsequently, we calculate the quantiles at significance levels  $\alpha$  and  $\beta$  for each of these sets. The quantiles corresponding to the scores in  $\mathbb{M}^g$  are denoted as  $q_\alpha^g$  and  $q_\beta^g$ , whereas those for the scores in  $\mathbb{M}^l$  are represented as  $q_\alpha^l$  and  $q_\beta^l$ . Ultimately, employing these computed quantiles, we establish two linear transformations,  $t^g(\cdot)$  and  $t^l(\cdot)$ , for the global and local anomaly maps, respectively. These transformations map the quantile at level  $\alpha$  to an anomaly score of 0 and at  $\beta$  to a score of 0.1, as

$$t^l(\mathbf{M}_i^l) = 0.1 \left( \mathbf{M}_i^l - \left( \frac{q_\alpha^l}{q_\beta^l - q_\alpha^l} \right) \mathbf{1}_{h^* \times w^*} \right), \quad \text{and} \quad t^g(\mathbf{M}_i^g) = 0.1 \left( \mathbf{M}_i^g - \left( \frac{q_\alpha^g}{q_\beta^g - q_\alpha^g} \right) \mathbf{1}_{h^* \times w^*} \right), \quad (9)$$

where  $\mathbf{1}_{h^* \times w^*}$  is a matrix of all ones. During inference, the global and local anomaly maps for each sample in  $\mathcal{D}_{\text{test}}$  are averaged after normalizing them with their respective transformations, which can be expressed as

$$\mathbf{M}_i = \frac{t^l(\mathbf{M}_i^l) + t^g(\mathbf{M}_i^g)}{2}.$$

For AD, the image anomaly score

$$s_i = \max_{h, w} \mathbf{M}_i(h, w) \quad (10)$$

is computed as the maximum value in the combined anomaly map  $\mathbf{M}_i$ . We follow [Batzner et al. \(2024\)](#) in mapping the destination values to 0 and 0.1 as they result in colour maps suitable for 0-to-1 colour scales. It should be noted that the choice of mapping the quantile values has no impact on the AU-ROC scores as it depends solely on the ranking of the scores.

## 4 Experimental Evaluation

In this section, we answer the following three questions: (i) How does ULSAD perform as compared to the SOTA methods? (ii) How effective is the local and global branch for the detection of *structural* and *logical* anomalies? (iii) How does each component in ULSAD impact the overall performance?

### 4.1 Setup

**Benchmark Datasets.** We evaluate our method on the following five IAD benchmarking datasets:

[1] **MVTec AD** ([Bergmann et al., 2019](#)). It consists of images from industrial manufacturing across 15 categories comprised of 10 objects and 5 textures. In totality, it contains 3,629 normal images for training. For evaluation, 1,258 anomalous images with varying pixel level defects and 467 normal images.

[2] **MVTec-LoCo** ([Bergmann et al., 2022](#)). An extension of MVTEC dataset, it encompasses both local structural anomalies and logical anomalies violating long-range dependencies. It consists of 5 categories, with 1,772 normal images for training and 304 normal images for validation. It also contains 1568 images, either normal or anomalous, for evaluation.

[3] **BTAD** ([Mishra et al., 2021](#)). It includes 3 categories with 1,799 normal images for training. It also

Table 1: Average Detection Performance in AUROC (%). Style: **best** and second best

Method	BTAD	MPDD	MVTec	MVTec-LOCO	VisA
PatchCore (Roth et al., 2022)	93.27	<u>93.27</u>	<b>98.75</b>	<u>81.49</u>	<u>91.48</u>
CFLOW (Gudovskiy et al., 2021)	93.57	87.11	94.47	73.62	87.77
DRÆM (Zavrtanik et al., 2021)	73.42	74.14	75.26	62.35	77.75
EfficientAD (Batzner et al., 2024)	88.26	85.42	<u>98.23</u>	80.62	91.21
FastFlow (Yu et al., 2021)	91.68	65.03	90.72	71.00	87.49
PaDiM (Defard et al., 2021)	93.20	68.48	91.25	68.38	83.28
Reverse Distillation (Deng & Li, 2022)	83.87	79.62	79.65	61.56	86.24
DFR (Yang et al., 2020)	<u>94.60</u>	79.75	93.54	72.87	85.18
<b>ULSAD (Ours)</b>	<b>96.17</b> $\pm$ 0.45	<b>95.73</b> $\pm$ 0.45	97.65 $\pm$ 0.38	<b>84.1</b> $\pm$ 0.86	<b>92.46</b> $\pm$ 0.45

consists of 290 anomalous images and 451 normal images for testing.

[4] **MPDD** (Jezek et al., 2021). It focuses on metal part fabrication defects. The images are captured in variable spatial orientation, position, and distance of multiple objects concerning the camera at different light intensities and with a non-homogeneous background. It consists of 6 classes of metal parts with 888 training images. For evaluation, the dataset has 176 normal and 282 anomalous images.

[5] **VisA** (Zou et al., 2022). It contains 10,821 high-resolution images (9,621 normal and 1,200 anomalous images) across 12 different categories. The anomalous images contain different types of anomalies such as scratches, bent, cracks, missing parts or misplacements. For each type of defect, there are 15-20 images, and an image can depict multiple defects.

**Evaluation metrics.** We measure the image-level anomaly detection performance via the area under the receiver operator curve (AUROC) based on the assigned anomaly score. To measure the anomaly localization performance, we use pixel-level AUROC and area under per region overlap curve (AUPRO). Furthermore, following prior works (Roth et al., 2022; Gudovskiy et al., 2021; Bergmann et al., 2019), we compute the average metrics over all the categories for each of the benchmark datasets. Moreover, for ULSAD, we report all the results over 5 runs with different random seeds.

**Baselines.** We compare our method with existing state-of-the-art unsupervised AD methods, namely PatchCore (Roth et al., 2022), PaDim (Defard et al., 2021), CFLOW (Gudovskiy et al., 2021), FastFLOW (Yu et al., 2021), DRÆM (Zavrtanik et al., 2021), Reverse Distillation (RD) (Deng & Li, 2022), EfficientAD (Batzner et al., 2024) and DFR (Yang et al., 2020). In this study, we only consider baselines that are capable of both anomaly detection and localization.

**Implementation details.** ULSAD is implemented in PyTorch (Paszke et al., 2019). For the baselines, we follow the implementation in Anomalib (Akçay et al., 2022), a widely used AD library for benchmarking. In ULSAD, we use a Wide-ResNet50-2 pre-trained on ImageNet (Zagoruyko & Komodakis, 2016) and **extract features from the second and third layers**, similar to PathCore (Roth et al., 2022). We use a convolutional autoencoder-like network in the global branch  $\phi$  and the feature prediction network  $\psi$ . It consists of convolution layers with LeakyReLU activation in the encoder and deconvolution layers in the decoder. The exact architecture is provided in the Appendix A. Unless otherwise stated, for all the experiments, we consider an image size of  $256 \times 256$ . We train ULSAD over 200 epochs for each category using an Adam optimizer with a learning rate of 0.0002 and a weight decay of 0.00002. For the baselines, we use the hyperparameters mentioned in the respective paper.

## 4.2 Evaluation Results

We summarize the anomaly detection performance of ULSAD in Table 1 and the localization performance in Table 2. On the BTAD dataset, we improve over the DFR by approximately 2% in detection. Inspecting the images from the dataset, we hypothesize that the difference stems from the use of a global branch in ULSAD as the structural imperfections are not limited to small regions. For localization, **Reverse Distillation** performs better owing to the use of multi-scale anomaly maps. We can observe similar improvements over

Table 2: Average Segmentation Performance in AUROC (%) and AUPRO (%). Style: **best** and second best

Method	BTAD	MPDD	MVTec	MVTec-LOCO	VisA
PatchCore (Roth et al., 2022)	96.85   71.48	<b>98.07</b>   90.84	<b>97.71</b>   91.15	75.77   69.09	97.93   85.12
CFLOW (Gudovskiy et al., 2021)	96.60   73.11	97.42   88.56	97.17   90.14	<u>76.99</u>   66.93	98.04   85.29
DR.EM (Zavrtanik et al., 2021)	59.04   22.48	86.96   70.04	75.01   49.72	63.69   40.06	71.31   54.68
EfficientAD (Batzner et al., 2024)	82.13   54.37	97.03   90.44	96.29   90.11	70.36   66.96	97.51   84.45
FastFlow (Yu et al., 2021)	96.15   75.27	93.60   76.89	96.44   88.79	75.55   53.04	97.32   81.70
PaDiM (Defard et al., 2021)	97.07   <u>77.80</u>	94.51   81.18	96.79   91.17	71.32   67.97	97.09   80.80
Reverse Distillation (Deng & Li, 2022)	<b>97.85</b>   <b>81.47</b>	<u>97.83</u>   <u>91.86</u>	97.25   <b>93.12</b>	68.55   66.28	<b>98.68</b>   <b>91.77</b>
DFR (Yang et al., 2020)	<u>97.62</u>   59.06	97.33   90.46	94.93   89.42	61.72   <u>69.78</u>	97.90   <u>91.72</u>
<b>ULSAD (Ours)</b>	96.73   75.41	97.45   <b>92.02</b>	<u>97.61</u>   <u>91.67</u>	<b>80.06</b>   <b>73.73</b>	<u>98.24</u>   87.12
	± 0.51   ± 3.95	± 0.99   ± 2.64	± 0.64   ± 1.36	± 0.20   ± 0.35	± 0.20   ± 0.89

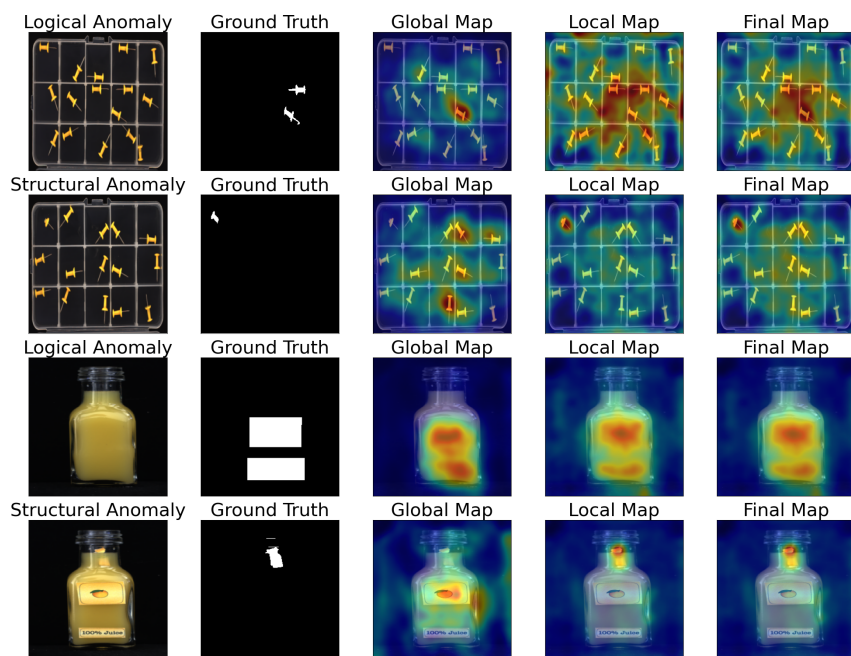


Figure 6: Example of anomaly maps obtained from global and local branches along with the combined map.

DFR on the MVTEC dataset. Although PatchCore provides superior performance on MVTEC, it should be noted that even without using a memory bank, ULSAD provides comparable results. Then, we focus on more challenging datasets such as MPDD and MVTEC-LOCO. While MPDD contains varying external conditions such as lighting, background and camera angles, MVTEC-LOCO contains both logical and structural anomalies. We can observe improvements over DFR ( $\sim 12 - 16\%$ ) in both datasets. This highlights the effectiveness of our method. We visualize the anomaly maps for samples from the “pushpin” and “juice bottle” categories in Figure 6. It can be seen that while the global branch is more suited to the detection of logical anomalies, the local branch is capable of detecting localized structural anomalies.

Overall, ULSAD demonstrates competitive results in anomaly detection compared to the baseline methods across all benchmark datasets. Additionally, the difference in performance between ULSAD and the baselines for anomaly localization is minimal. The most notable difference is in the AUPRO score on the BTAD, MVTEC, and VisA datasets. Nonetheless, while the SOTA method provides slightly better performance in the localization of structural anomalies, ULSAD provides similar performance across both logical and structural anomalies.

We present anomaly maps obtained from different methods in Figure 9. Extended versions of Tables 1 and 2 are provided in Appendix B. Additionally, we provide the results on MVTECLOCO, split between logical and structural anomalies, in Appendix B.1.

### 4.3 Ablation study

**Analysis of main components.** We investigate the impact of the key components in ULSAD by sequentially adding them. For the ablations, we report the detection and localization performance on the MVTECLOCO dataset in Table 3. In the initial set of experiments, we set both  $\lambda_l$  and  $\lambda_g$  to 0. This results in the use of only the difference in values while computing the global and local objectives. The first row of the table corresponds to the use of only the local branch. We observe a slight improvement in both detection and localization performance when considering the global branch in addition to the local branch. For the sake of completeness, we also consider here a variant of the global loss  $\mathcal{L}_{pg}$  where we compute the  $\ell_2$  distance between the feature maps instead of computing the self- and cross-attention maps. We refer to the alternative in the table as  $\mathcal{L}_{pg}^d$ . It can be seen that when using the cross-attention map without considering the difference in direction, the performance drops compared to using  $\mathcal{L}_{pg}^d$ .

The global loss  $\mathcal{L}_{pg}$  is incorporated to learn the spatial relationships among objects in normal images for detecting logical anomalies. However, when used in isolation, it limits ULSAD’s performance to detecting only logical anomalies and fails to capture localized structural anomalies. Additionally, as discussed in Section 3.3, the global branch is prone to false positives in the presence of sharp edges or heavily textured surfaces. Therefore, the incorporation of the loss term  $\mathcal{L}_{lg}$ , which relates to the global and local branches, results in a significant boost in overall performance.

Further, the difference in performance between the two variants of the global loss becomes negligible. Considering the difference of direction in both the global and local objectives, we observe improvements over just using the difference in values, as seen in the last five rows of the table. Overall, the best performance in detection is obtained by incorporating both the local-global loss term and the attention-based global objective.

**Effect of backbone.** We investigate the impact of using different pre-trained backbones in ULSAD in Figure 7. It is observed that the overall best performance is obtained by using a Wide-ResNet101-2 architecture in both detection and localization. More specifically, for detection, Wide-ResNet variants are more effective than the ResNet architectures, whereas, for localization performance measured using Pixel AUROC, the deeper networks such as ResNet152 and Wide-ResNet101-2 seem to have precedence over their shallower counterparts. Overall, we can see that performance is robust to the choice of pre-trained model architecture. In our experiments, we utilize a Wide-ResNet50-2 architecture which is used by most of our baselines for fair comparison.

**Effect of normalization.** We analyze the impact of the quantile-based normalization on the performance metrics by considering multiple values for  $\alpha$  and  $\beta$ . The results are shown in the Figure 8. It can be seen that the final performance is robust to the choice of  $\alpha$  and  $\beta$ . For our experiments, we set  $\alpha = 0.9$  and  $\beta = 0.995$ .

Table 3: Ablation of the main components of ULSAD.

Local Branch $\lambda_l$	Global Branch				Performance (%)		
	$\lambda_g$	$\mathcal{L}_{lg}$	$\mathcal{L}_{pg}^d$	$\mathcal{L}_{pg}$	IAUROC	PAUROC	PAUPRO
0.0	-	-	-	-	77.67	75.17	73.37
0.0	0.0	-	✓	-	77.69	79.77	75.26
0.0	0.0	-	-	✓	71.67	73.92	67.22
0.0	0.0	✓	✓	-	81.40	82.12	77.47
0.0	0.0	✓	-	✓	81.08	81.97	76.45
0.5	-	-	-	-	79.14	76.57	73.41
0.5	0.5	-	✓	-	80.50	81.85	77.35
0.5	0.5	-	-	✓	74.51	76.59	69.01
0.5	0.5	✓	✓	-	82.19	81.25	75.50
0.5	0.5	✓	-	✓	<b>84.10</b>	<b>80.06</b>	<b>73.73</b>
					$\pm 0.86$	$\pm 0.20$	$\pm 0.35$

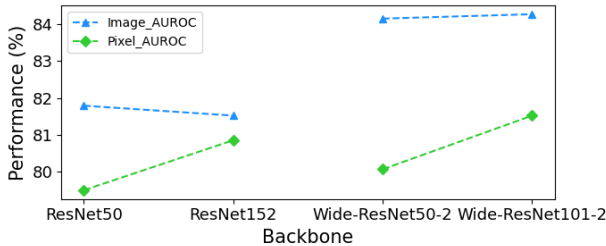


Figure 7: Ablation study of the backbone network

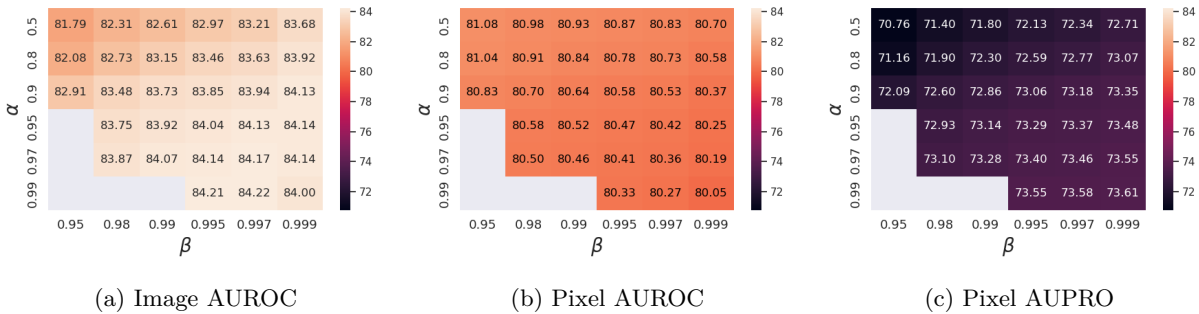
Figure 8: Ablation of  $\alpha$  and  $\beta$ .

Table 4: Memory and Computational Efficiency on MVTecLOCO dataset.

	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
I-AUROC $\uparrow$	73.62	62.35	71.00	68.38	<u>81.49</u>	61.56	72.87	80.62	<b>84.10</b> $\pm$ 0.86
P-AUROC $\uparrow$	<u>76.99</u>	63.69	75.55	71.32	75.77	68.55	61.72	70.36	<b>80.06</b> $\pm$ 0.20
P-AUPRO $\uparrow$	66.93	40.06	53.04	67.97	69.09	66.28	<u>69.78</u>	66.96	<b>73.73</b> $\pm$ 0.35
Throughput (img / s) $\uparrow$	11.69	10.06	30.21	<u>33.45</u>	32.70	<b>34.87</b>	15.08	23.33	33.42
GPU Memory (GB) $\downarrow$	2.57	7.95	<b>1.69</b>	<u>1.92</u>	6.80	1.93	5.85	3.48	2.17

## 5 Memory and Computational Complexity

We report the computational cost and memory requirements of ULSAD compared to the baselines in Table 4. For this analysis, we ran inference on the test samples in the MVTecLOCO dataset using an NVIDIA A100 GPU. We measured throughput with a batch size of 32, as a measure of computational complexity, following EfficientAD (Batzner et al., 2024). Throughput is defined as the number of images processed per second when processing in batches. ULSAD demonstrates higher throughput than most baselines while maintaining competitive anomaly detection and localization performance.

In addition to throughput, we also report peak GPU memory usage in Table 4 to highlight the memory efficiency of ULSAD. It is evident that ULSAD requires approximately one-third of the memory compared to retrieval-based methods such as PatchCore, which is one of the state-of-the-art methods for IAD.

For DFR (Yang et al., 2020), we followed the authors’ approach by using a multiscale representation, concatenating features from 12 layers of the pre-trained network  $\theta$  for anomaly detection. This approach results in reduced throughput and increased memory usage, as shown in Table 4. With our proposed modifications in ULSAD, we achieve superior performance using features from only 2 layers of  $\theta$ , drastically reducing memory requirements to approximately one-third of DFR’s and increasing throughput by approximately two times.

## 6 Limitations

For training ULSAD, we assume, similar to prior works in unsupervised anomaly detection (Ruff et al., 2021; Chandola et al., 2009; Roth et al., 2022; Batzner et al., 2024), that the training dataset is "clean", i.e., it does not contain any anomalous samples. This setting is referred to in the literature as one-class classification (Ruff et al., 2018). However, since anomalies are unknown a priori, such an assumption might affect performance in real-world applications. The study of the effect of contamination (Wang et al., 2019; Jiang et al., 2022; Yoon et al., 2022) is an active area of research and is beyond the scope of our current study. We leave for future research the analysis of contamination’s impact on ULSAD’s performance and the exploration of options to make the learning process robust in the presence of anomalies.

## 7 Conclusion

Our study focuses on Deep Feature Reconstruction (DFR), which offers a memory- and compute-efficient way of detecting structural anomalies. We develop ULSAD, a unified framework for detecting structural and logical anomalies, which utilizes a dual-branch architecture by extending DFR. Specifically, we improve the

training objective of DFR to enhance the capability to detect structural anomalies, which we refer to as the local branch. Furthermore, the bottleneck used in FRN at the local branch acts as an implicit regularizer, preventing the over-expressiveness of the model without using additional computationally expensive methods of regularization. Additionally, we introduce an attention-based loss using an autoencoder-like network for detecting logical anomalies in the global branch. Through extensive experimentation across five benchmarking image AD datasets, we demonstrate that **ULSAD** can provide competitive performance in the detection and localization of both structural and logical anomalies compared to eight state-of-the-art methods. Importantly, it performs well even when compared to memory-intensive retrieval-based methods such as PatchCore (Roth et al., 2022). Lastly, ablations demonstrate the impact of various components used in **ULSAD** and the influence of the pre-trained backbone  $\theta$  on the overall performance.

**Reproducibility Statement.** We provide extensive descriptions of implementation details (in Section 4.1) and algorithm (in Algorithm 1) to help readers reproduce our results. Every measure is taken to ensure fairness in our comparisons by adopting the most commonly adopted evaluation settings in the anomaly detection literature. More specifically, we use the [Anomalib library](#) for our experiments, whenever applicable, for comparing their performance with ULSAD. For methods such as DFR (Yang et al., 2020), which is not implemented in Anomalib, we refer to the original codebase provided by the authors. We will publicly release the code upon acceptance.

**Ethics Statement.** We have read the TMLR Ethics Guidelines (<https://jmlr.org/tmlr/ethics.html>) and ensured that this work adheres to it. All benchmark datasets and pre-trained model checkpoints are publicly available and not directly subject to ethical concerns.

## References

- Charu C Aggarwal, Alexander Hinneburg, and Daniel A Keim. On the surprising behavior of distance metrics in high dimensional space. In *Database Theory—ICDT 2001: 8th International Conference London, UK, January 4–6, 2001 Proceedings 8*, pp. 420–434. Springer, 2001.
- Mohiuddin Ahmed, Abdun Naser Mahmood, and Md Rafiqul Islam. A survey of anomaly detection techniques in financial domain. *Future Generation Computer Systems*, 55:278–288, 2 2016. ISSN 0167-739X. doi: 10.1016/J.FUTURE.2015.01.001.
- Samet Akcay, Dick Ameln, Ashwin Vaidya, Barath Lakshmanan, Nilesh Ahuja, and Utku Genc. Anomalib: A Deep Learning Library for Anomaly Detection. In *2022 IEEE International Conference on Image Processing (ICIP)*, pp. 1706–1710. IEEE, 10 2022. ISBN 978-1-6654-9620-9. doi: 10.1109/ICIP46576.2022.9897283. URL <https://ieeexplore.ieee.org/document/9897283/>.
- Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael Rabbat, Yann LeCun, and Nicolas Ballas. Self-supervised learning from images with a joint-embedding predictive architecture. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15619–15629, 2023.
- Kilian Batzner, Lars Heckler, and Rebecca König. EfficientAD: Accurate Visual Anomaly Detection at Millisecond-Level Latencies. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pp. 128–138, 1 2024. URL <https://arxiv.org/abs/2303.14535v2>.
- Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. MVTEC ad-A comprehensive real-world dataset for unsupervised anomaly detection. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 9584–9592, 2019. ISSN 10636919. doi: 10.1109/CVPR.2019.00982.
- Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. Uninformed Students: Student-Teacher Anomaly Detection With Discriminative Latent Embeddings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4183–4192, 2020. URL [www.mvtec.com](http://www.mvtec.com).
- Paul Bergmann, Kilian Batzner, Michael Fauser, David Sattlegger, and Carsten Steger. Beyond Dents and Scratches: Logical Constraints in Unsupervised Anomaly Detection and Localization. *International Journal of Computer Vision*, 130(4):947–969, 2022. ISSN 1573-1405. doi: 10.1007/s11263-022-01578-9. URL <https://doi.org/10.1007/s11263-022-01578-9>.
- Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection. *ACM Computing Surveys*, 41(3): 1–58, 7 2009. ISSN 0360-0300. doi: 10.1145/1541880.1541882. URL <https://dl.acm.org/doi/10.1145/1541880.1541882>.
- Heejeong Choi, Donghwa Kim, Jounghee Kim, Jina Kim, and Pilsung Kang. Explainable anomaly detection framework for predictive maintenance in manufacturing systems. *Applied Soft Computing*, 125:109147, 8 2022. ISSN 1568-4946. doi: 10.1016/J.ASOC.2022.109147.

- Thomas Defard, Aleksandr Setkov, Angélique Loesch, and Romaric Audigier. PaDiM: a Patch Distribution Modeling Framework for Anomaly Detection and Localization. In *International Conference on Pattern Recognition*, volume 12664 LNCS, pp. 475–489. Springer Science and Business Media Deutschland GmbH, 11 2021. ISBN 9783030687984. doi: 10.48550/arxiv.2011.08785. URL <https://arxiv.org/abs/2011.08785v1>.
- Hanqiu Deng and Xingyu Li. Anomaly Detection via Reverse Distillation From One-Class Embedding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9737–9746, 2022.
- Alexey Dosovitskiy and Thomas Brox. Generating images with perceptual similarity metrics based on deep networks. *Advances in neural information processing systems*, 29, 2016.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- Tharindu Fernando, Harshala Gammulle, Simon Denman, Sridha Sridharan, and Clinton Fookes. Deep learning for medical anomaly detection – a survey. *ACM Comput. Surv.*, 54(7), jul 2021. ISSN 0360-0300. doi: 10.1145/3464423. URL <https://doi.org/10.1145/3464423>.
- Han Gao, Huiyuan Luo, Fei Shen, and Zhengtao Zhang. Towards Total Online Unsupervised Anomaly Detection and Localization in Industrial Vision, 5 2023. URL <https://arxiv.org/abs/2305.15652v1>.
- Dong Gong, Lingqiao Liu, Vuong Le, Budhaditya Saha, Moussa Reda Mansour, Svetha Venkatesh, and Anton Van Den Hengel. Memorizing Normality to Detect Anomaly: Memory-augmented Deep Autoencoder for Unsupervised Anomaly Detection. *Proceedings of the IEEE International Conference on Computer Vision*, 2019-October:1705–1714, 4 2019. ISSN 15505499. doi: 10.48550/arxiv.1904.02639. URL <https://arxiv.org/abs/1904.02639v2>.
- Denis Gudovskiy, Shun Ishizaka, and Kazuki Kozuka. CFLOW-AD: Real-Time Unsupervised Anomaly Detection with Localization via Conditional Normalizing Flows. *Proceedings - 2022 IEEE/CVF Winter Conference on Applications of Computer Vision, WACV 2022*, pp. 1819–1828, 7 2021. doi: 10.1109/WACV51458.2022.00188. URL <https://arxiv.org/abs/2107.12571v1>.
- Waleed Hilal, S. Andrew Gadsden, and John Yawney. Financial Fraud: A Review of Anomaly Detection Techniques and Recent Advances. *Expert Systems with Applications*, 193:116429, 5 2022. ISSN 0957-4174. doi: 10.1016/J.ESWA.2021.116429.
- Stepan Jezek, Martin Jonak, Radim Burget, Pavel Dvorak, and Milos Skotak. Deep learning-based defect detection of metal parts: Evaluating current methods in complex conditions. *International Congress on Ultra Modern Telecommunications and Control Systems and Workshops*, 2021-October:66–71, 2021. ISSN 2157023X. doi: 10.1109/ICUMT54235.2021.9631567.
- Xi Jiang, Jianlin Liu, Jinbao Wang, Qiang Nie, Kai WU, Yong Liu, Chengjie Wang, and Feng Zheng. SoftPatch: Unsupervised Anomaly Detection with Noisy Data. *Advances in Neural Information Processing Systems*, 35: 15433–15445, 2022. URL <https://github.com/TencentYoutuResearch/AnomalyDetection-SoftPatch>.
- Chieh-Hsin Lai, Dongmian Zou, and Gilad Lerman. Robust Subspace Recovery Layer for Unsupervised Anomaly Detection. *arXiv preprint*, 3 2019. doi: 10.48550/arxiv.1904.00152. URL <https://arxiv.org/abs/1904.00152v2>.
- Sungwook Lee, Seunghyun Lee, and Byung Cheol Song. CFA: Coupled-hypersphere-based Feature Adaptation for Target-Oriented Anomaly Localization. *IEEE Access*, 10:78446–78454, 6 2022. ISSN 21693536. doi: 10.1109/ACCESS.2022.3193699. URL <https://arxiv.org/abs/2206.04325v1>.
- Jiaqi Liu, Guoyang Xie, Jinbao Wang, Shangnian Li, Chengjie Wang, Feng Zheng, and Yaochu Jin. Deep industrial image anomaly detection: A survey. *Machine Intelligence Research*, 21(1):104–135, 2024.



- Pankaj Mishra, Riccardo Verk, Daniele Fornasier, Claudio Piciarelli, and Gian Luca Foresti. VT-ADL: A Vision Transformer Network for Image Anomaly Detection and Localization. *IEEE International Symposium on Industrial Electronics*, 2021-June, 4 2021. doi: 10.1109/ISIE45552.2021.9576231. URL <http://arxiv.org/abs/2104.10036><http://dx.doi.org/10.1109/ISIE45552.2021.9576231>.
- Duc Tam Nguyen, Zhongyu Lou, Michael Klar, and Thomas Brox. Anomaly Detection With Multiple-Hypotheses Predictions. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 4800–4809. PMLR, 7 2019. URL <https://proceedings.mlr.press/v97/nguyen19b.html>.
- Guansong Pang, Chunhua Shen, Longbing Cao, and Anton van den Hengel. Deep Learning for Anomaly Detection: A Review. *ACM Computing Surveys*, 54(2), 7 2020. doi: 10.1145/3439950. URL <http://arxiv.org/abs/2007.02500><http://dx.doi.org/10.1145/3439950>.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury Google, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf Xamla, Edward Yang, Zach Devito, Martin Raison Nabla, Alykhan Tejani, Sasank Chilamkurthy, Qure Ai, Benoit Steiner, Lu Fang Facebook, Junjie Bai Facebook, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning Library. *Advances in Neural Information Processing Systems*, 32, 2019.
- Pramuditha Perera, Ramesh Nallapati, and Bing Xiang. OCGAN: One-class novelty detection using gans with constrained latent representations. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2019-June:2893–2901, 6 2019. ISSN 10636919. doi: 10.1109/CVPR.2019.00301.
- Tal Reiss, Niv Cohen, Eliahu Horwitz, Ron Abutbul, and Yedid Hoshen. Anomaly Detection Requires Better Representations. In *European Conference on Computer Vision*, pp. 56–68, 10 2022. ISBN 9783031250682. doi: 10.1007/978-3-031-25069-9{\\\_}4. URL <https://arxiv.org/abs/2210.10773v1>.
- Karsten Roth, Latha Pemula, Joaquin Zepeda, Bernhard Schölkopf, Thomas Brox, and Peter Gehler. Towards Total Recall in Industrial Anomaly Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14318–14328, 2022. doi: 10.48550/arxiv.2106.08265. URL <https://arxiv.org/abs/2106.08265v2>.
- Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. Deep One-Class Classification. In *International Conference on Machine Learning*, pp. 4393–4402. PMLR, 2018. URL <https://proceedings.mlr.press/v80/ruff18a.html>.
- Lukas Ruff, Jacob R. Kauffmann, Robert A. Vandermeulen, Gregoire Montavon, Wojciech Samek, Marius Kloft, Thomas G. Dietterich, and Klaus Robert Muller. A Unifying Review of Deep and Shallow Anomaly Detection. *Proceedings of the IEEE*, 109(5):756–795, 5 2021. ISSN 15582256. doi: 10.1109/JPROC.2021.3052449.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3):211–252, 12 2015. ISSN 15731405. doi: 10.1007/S11263-015-0816-Y/FIGURES/16. URL <https://link.springer.com/article/10.1007/s11263-015-0816-y>.
- Mohammad Sabokrou, Mohammad Khalooei, Mahmood Fathy, and Ehsan Adeli. Adversarially Learned One-Class Classifier for Novelty Detection. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 3379–3388, 12 2018. ISSN 10636919. doi: 10.1109/CVPR.2018.00356.
- Mohammadreza Salehi, Niousha Sadjadi, Soroosh Baselizadeh, Mohammad H. Rohban, and Hamid R. Rabiee. Multiresolution knowledge distillation for anomaly detection. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 14897–14907, 11 2021. ISSN 10636919. doi: 10.1109/CVPR46437.2021.01466. URL <https://arxiv.org/abs/2011.11108v1>.

- Mohammadreza Salehi, Hossein Mirzaei, Dan Hendrycks, Yixuan Li, Mohammad Hossein Rohban, and Mohammad Sabokrou. A unified survey on anomaly, novelty, open-set, and out of-distribution detection: Solutions and future challenges. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL <https://openreview.net/forum?id=aRtjVZvbpK>.
- Wuqin Tang, Qiang Yang, Kuixiang Xiong, and Wenjun Yan. Deep learning based automatic defect identification of photovoltaic module using electroluminescence images. *Solar Energy*, 201:453–460, 5 2020. ISSN 0038-092X. doi: 10.1016/J.SOLENER.2020.03.049.
- Robert Tibshirani and Trevor Hastie. Outlier sums for differential gene expression analysis. *Biostatistics*, 8(1):2–8, 1 2007. ISSN 1465-4644. doi: 10.1093/BIostatistics/KXL005. URL <https://academic.oup.com/biostatistics/article/8/1/2/252461>.
- Frederick Tung and Greg Mori. Similarity-preserving knowledge distillation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 1365–1374, 2019.
- Siqi Wang, Yijie Zeng, Xinwang Liu, En Zhu, Jianping Yin, Chuanfu Xu, and Marius Kloft. Effective End-to-end Unsupervised Outlier Detection via Inlier Priority of Discriminative Network. *Advances in Neural Information Processing Systems*, 32, 2019.
- Xudong Yan, Huaidong Zhang, Xuemiao Xu, Xiaowei Hu, and Pheng-Ann Heng. Learning Semantic Context from Normal Samples for Unsupervised Anomaly Detection. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(4):3110–3118, 5 2021. URL <https://ojs.aaai.org/index.php/AAAI/article/view/16420>.
- Jie Yang, Yong Shi, and Zhiqian Qi. DFR: Deep Feature Reconstruction for Unsupervised Anomaly Segmentation. *Neurocomputing*, 424:9–22, 12 2020. doi: 10.1016/j.neucom.2020.11.018. URL <http://arxiv.org/abs/2012.07122><http://dx.doi.org/10.1016/j.neucom.2020.11.018>.
- Jinsung Yoon, Kihyuk Sohn, Chun-Liang Li, Sercan O. Arik, Chen-Yu Lee, and Tomas Pfister. Self-supervise, Refine, Repeat: Improving Unsupervised Anomaly Detection. *Transactions on Machine Learning Research (TMLR)*, 2022. URL <https://arxiv.org/abs/2106.06115v2>.
- Jinsung Yoon, Kihyuk Sohn, Chun-Liang Li, Sercan O. Arik, and Tomas Pfister. SPADE: Semi-supervised Anomaly Detection under Distribution Mismatch. *Transactions on Machine Learning Research*, 11 2023. URL <http://arxiv.org/abs/2212.00173>.
- Jiawei Yu, Ye Zheng, Xiang Wang, Wei Li, Yushuang Wu, Rui Zhao, and Liwei Wu. FastFlow: Unsupervised Anomaly Detection and Localization via 2D Normalizing Flows. *arXiv preprint arXiv:2111.07677*, 11 2021. doi: 10.48550/arxiv.2111.07677. URL <https://arxiv.org/abs/2111.07677v2>.
- Sergey Zagoruyko and Nikos Komodakis. Wide Residual Networks. *British Machine Vision Conference 2016, BMVC 2016*, 2016-September:1–87, 5 2016. doi: 10.5244/C.30.87. URL <https://arxiv.org/abs/1605.07146v4>.
- Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *International Conference on Learning Representations*, 2017. URL [https://openreview.net/forum?id=Sks9\\_ajex](https://openreview.net/forum?id=Sks9_ajex).
- Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. DRÆM - A discriminatively trained reconstruction embedding for surface anomaly detection. *Proceedings of the IEEE International Conference on Computer Vision*, pp. 8310–8319, 8 2021. ISSN 15505499. doi: 10.1109/ICCV48922.2021.00822. URL <https://arxiv.org/abs/2108.07610v2>.
- Houssam Zenati, Chuan Sheng Foo, Bruno Lecouat, Gaurav Manek, and Vijay Ramaseshan Chandrasekhar. Efficient GAN-Based Anomaly Detection. *arXiv preprint*, 2 2018. doi: 10.48550/arxiv.1802.06222. URL <https://arxiv.org/abs/1802.06222v2>.

Jie Zhang, Masanori Suganuma, and Takayuki Okatani. Contextual Affinity Distillation for Image Anomaly Detection. *arXiv preprint arXiv:2307.03101*, 7 2023. URL <https://arxiv.org/abs/2307.03101v1>.

Kang Zhou, Yuting Xiao, Jianlong Yang, Jun Cheng, Wen Liu, Weixin Luo, Zaiwang Gu, Jiang Liu, and Shenghua Gao. Encoding Structure-Texture Relation with P-Net for Anomaly Detection in Retinal Images. In *European conference on computer vision*, volume 12365 LNCS, pp. 360–377. Springer, 2020. doi: 10.1007/978-3-030-58565-5\_{\\_}22. URL [https://link.springer.com/10.1007/978-3-030-58565-5\\_22](https://link.springer.com/10.1007/978-3-030-58565-5_22).

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

## A Implementation Details

ULSAD is implemented in PyTorch (Paszke et al., 2019). Specifically, we used the Anomalib (Akçay et al., 2022) library by incorporating our code within it. It helps us have a fair comparison as we use the implementations of baselines from Anomalib. Moreover, we used a single NVIDIA A4000 GPU for all the experiments **unless mentioned otherwise**. The architecture of FRN and global autoencoder-like model is provided in Table 5 and 6, respectively.

Table 5: Feature Reconstruction Network of ULSAD.

	Layer Name	Stride	Kernel Size	Number of Kernels	Padding	Activation
Encoder	Conv-1	2	$3 \times 3$	768	1	ReLU
	BatchNorm-1	-	-	-	-	-
	Conv-2	2	$3 \times 3$	1536	1	ReLU
	BatchNorm-2	-	-	-	-	-
	Conv-3	1	$3 \times 3$	1536	1	ReLU
	BatchNorm-3	-	-	-	-	-
Decoder	ConvTranspose-1	2	$4 \times 4$	768	1	ReLU
	BatchNorm-4	-	-	-	-	-
	ConvTranspose-2	2	$4 \times 4$	384	1	ReLU
	BatchNorm-5	-	-	-	-	-
	ConvTranspose-3	1	$5 \times 5$	384	1	ReLU
	BatchNorm-6	-	-	-	-	-

Table 6: Global Autoencoder of ULSAD.

	Layer Name	Stride	Kernel Size	Number of Kernels	Padding	Activation
Encoder	Conv-1	2	$4 \times 4$	32	1	ReLU
	BatchNorm-1	-	-	-	-	-
	Conv-2	2	$4 \times 4$	32	1	ReLU
	BatchNorm-2	-	-	-	-	-
	Conv-3	2	$4 \times 4$	64	1	ReLU
	BatchNorm-3	-	-	-	-	-
	Conv-4	2	$4 \times 4$	64	1	ReLU
	BatchNorm-4	-	-	-	-	-
	Conv-5	2	$4 \times 4$	64	1	ReLU
	BatchNorm-5	-	-	-	-	-
	Conv-6	1	$8 \times 8$	64	1	ReLU
	BatchNorm-6	-	-	-	-	-
Decoder	Interpolate-1 (31, mode= "bilinear")	-	-	-	-	-
	Conv-1	1	$4 \times 4$	64	2	ReLU
	BatchNorm-1	-	-	-	-	-
	Interpolate-2 (8, mode= "bilinear")	-	-	-	-	-
	Conv-2	1	$4 \times 4$	64	2	ReLU
	BatchNorm-2	-	-	-	-	-
	Interpolate-3 (16, mode= "bilinear")	-	-	-	-	-
	Conv-3	1	$4 \times 4$	64	2	ReLU
	BatchNorm-3	-	-	-	-	-
	Interpolate-4 (32, mode= "bilinear")	-	-	-	-	-
	Conv-4	1	$4 \times 4$	64	2	ReLU
	BatchNorm-4	-	-	-	-	-
	Interpolate-5 (64, mode= "bilinear")	-	-	-	-	-
	Conv-5	1	$4 \times 4$	64	2	ReLU
	BatchNorm-5	-	-	-	-	-
	Interpolate-6 (32, mode= "bilinear")	-	-	-	-	-
	Conv-6	1	$3 \times 3$	64	1	ReLU
BatchNorm-6	-	-	-	-	-	
Conv-7	1	$3 \times 3$	384	1	ReLU	

## B Extended Results

Extended versions of the Table 1 and 2 are provided in Tables 7-21. It shows the performance of ULSAD per category of the benchmark datasets for anomaly detection and localization. Additionally, we provide a visual comparison of the generated anomaly maps using the MVTecLOCO dataset in Figure 9.

Table 7: Anomaly detection based on Image AUROC on MVTec dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
bottle	<b>100.0</b>	94.6	98.57	<u>99.37</u>	<b>100.0</b>	98.10	94.92	<b>100.0</b>	<b>100.0</b> ± 0.00
cable	93.46	74.29	89.30	86.94	<b>98.54</b>	95.41	79.54	94.61	<u>97.92</u> ± 0.18
capsule	91.74	65.46	86.04	88.43	<b>97.93</b>	81.01	<u>96.01</u>	95.57	94.61 ± 0.28
carpet	93.26	57.70	<u>98.76</u>	97.31	97.91	57.95	97.59	<b>99.70</b>	98.50 ± 0.17
grid	93.57	76.02	<u>99.08</u>	84.04	97.24	93.07	94.57	<b>100.0</b>	92.67 ± 1.20
hazelnut	<b>100.0</b>	84.43	81.57	86.07	<b>100.0</b>	99.82	<b>100.0</b>	93.39	<u>99.93</u> ± 0.05
leather	<u>99.97</u>	79.31	<b>100.0</b>	99.66	<b>100.0</b>	42.39	99.46	99.92	<b>100.0</b> ± 0.00
metal_nut	<b>99.76</b>	45.06	94.53	96.92	<u>99.61</u>	67.20	93.06	99.34	98.88 ± 0.07
pill	90.73	44.65	87.53	88.52	94.35	54.66	92.06	<b>99.08</b>	<u>96.17</u> ± 0.39
screw	88.30	89.38	66.00	75.24	<b>98.26</b>	94.57	93.69	<u>98.09</u>	95.20 ± 0.15
tile	<b>100.0</b>	90.15	95.42	95.49	98.67	97.37	92.97	99.85	<u>99.99</u> ± 0.02
toothbrush	83.33	80.28	79.44	<u>93.61</u>	<b>100.0</b>	84.72	<b>100.0</b>	<b>100.0</b>	<b>100.0</b> ± 0.00
transistor	91.50	88.37	94.42	92.29	<b>100.0</b>	83.29	80.54	96.57	<u>97.65</u> ± 0.57
wood	98.33	90.96	97.54	98.33	<b>99.30</b>	53.16	98.77	98.13	<u>98.81</u> ± 0.23
zipper	93.07	68.17	92.54	86.48	<b>99.47</b>	92.04	89.97	<u>99.22</u>	94.36 ± 0.13
Mean	94.47	75.26	90.72	91.25	<b>98.75</b>	79.65	93.54	<u>98.23</u>	97.65 ± 0.38

Table 8: Anomaly segmentation performance based on Pixel AUROC on MVTec dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
bottle	<b>98.58</b>	76.53	97.8	98.3	97.98	<u>98.31</u>	90.83	<u>98.31</u>	96.21 ± 2.21
cable	96.1	66.59	95.71	96.81	<u>98.03</u>	96.37	91.37	<b>98.5</b>	97.71 ± 0.06
capsule	98.71	86.96	98.37	98.67	98.77	<b>98.96</b>	98.46	98.33	<u>98.95</u> ± 0.03
carpet	98.57	71.95	98.27	98.68	98.67	<u>99.05</u>	98.46	94.83	<b>99.18</b> ± 0.06
grid	97.49	53.56	<u>98.39</u>	92.82	97.86	<b>99.01</b>	97.41	96.02	95.47 ± 1.09
hazelnut	98.64	84.66	94.79	97.85	98.43	<b>98.91</b>	98.53	96.15	<u>98.81</u> ± 0.03
leather	<u>99.42</u>	63.32	<b>99.62</b>	99.30	98.87	99.17	99.33	97.5	98.68 ± 0.01
metal_nut	97.97	80.25	97.01	96.71	<b>98.51</b>	97.68	93.02	<u>98.07</u>	97.62 ± 0.03
pill	<u>97.83</u>	77.17	96.38	95.03	97.53	96.96	96.86	<b>98.63</b>	96.67 ± 0.09
screw	97.64	83.38	89.87	97.89	99.19	<b>99.43</b>	99.07	98.50	<u>99.33</u> ± 0.01
tile	<b>96.68</b>	85.75	93.14	92.42	94.86	95.47	90.82	91.61	<u>95.78</u> ± 0.05
toothbrush	98.16	90.70	97.50	<u>98.83</u>	98.67	<b>98.99</b>	98.49	96.0	98.42 ± 0.02
transistor	89.91	63.23	96.45	<u>96.85</u>	96.84	86.77	79.11	94.77	<b>98.89</b> ± 0.05
wood	94.70	71.73	<b>95.71</b>	93.83	93.31	95.06	<u>95.36</u>	90.85	95.20 ± 0.31
zipper	97.08	69.31	97.62	97.82	<u>98.06</u>	<b>98.54</b>	96.85	96.21	97.24 ± 0.07
Mean	97.17	75.01	96.44	96.79	<b>97.71</b>	97.25	94.93	96.29	<u>97.61</u> ± 0.64

Table 9: Anomaly segmentation performance based on Pixel AUPRO on MVTec dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
bottle	94.19	50.05	92.0	<u>95.11</u>	92.28	<b>95.12</b>	83.14	93.84	90.16 ± 3.24
cable	85.85	28.58	86.65	89.65	<u>90.77</u>	90.32	83.09	<b>92.53</b>	88.63 ± 0.48
capsule	90.47	81.11	90.15	92.62	92.4	<u>93.93</u>	<b>96.33</b>	91.09	93.77 ± 0.16
carpet	92.64	48.64	94.63	95.59	92.7	<b>96.41</b>	95.47	90.99	<u>96.39</u> ± 0.24
grid	90.52	17.71	<u>93.95</u>	82.52	89.46	<b>96.39</b>	91.15	93.14	83.3 ± 3.96
hazelnut	96.12	76.19	93.92	92.95	94.44	<u>96.92</u>	<b>97.17</b>	83.25	94.87 ± 0.3
leather	<u>98.39</u>	52.1	<b>99.06</b>	97.91	96.33	97.97	98.34	97.32	97.44 ± 0.01
metal_nut	88.97	35.79	85.89	90.45	91.9	<b>94.4</b>	87.01	<u>92.97</u>	91.58 ± 0.19
pill	93.67	64.26	91.0	93.88	93.92	94.76	<u>95.86</u>	<b>95.93</b>	94.5 ± 0.08
screw	90.25	53.22	68.6	92.14	95.39	<b>97.05</b>	95.96	96.04	<u>96.45</u> ± 0.1
tile	<b>91.49</b>	58.48	81.01	78.32	79.64	<u>88.4</u>	79.36	83.54	87.82 ± 0.15
toothbrush	81.05	54.02	80.62	<b>93.52</b>	86.48	92.23	<u>92.93</u>	88.61	86.28 ± 0.46
transistor	78.75	51.37	88.92	89.04	<b>94.06</b>	75.05	64.25	82.82	<u>91.61</u> ± 0.68
wood	90.5	45.29	<b>93.26</b>	91.39	85.08	<u>92.69</u>	92.48	76.16	91.34 ± 0.29
zipper	89.3	28.98	92.12	92.48	92.43	<b>95.18</b>	88.74	<u>93.48</u>	90.89 ± 0.35
Mean	90.14	49.72	88.79	91.17	91.15	<b>93.12</b>	89.42	90.11	<u>91.67</u> ± 1.36

Table 10: Anomaly detection performance based on Image AUROC on MVTecLOCO dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
breakfast_box	71.86	70.26	74.04	63.66	<b>85.24</b>	52.69	65.46	74.80	<u>83.54</u> ± 0.23
juice_bottle	81.70	62.55	78.03	88.42	<u>92.51</u>	76.28	86.81	98.89	<b>97.12</b> ± 0.10
pushpins	73.43	51.32	61.20	61.30	75.54	50.72	72.68	<u>80.58</u>	<b>86.85</b> ± 0.94
screw_bag	65.48	59.39	68.04	60.14	<u>69.90</u>	65.15	63.55	67.42	<b>70.71</b> ± 1.49
splicing_connectors	75.63	68.25	73.71	68.40	<b>84.24</b>	62.95	75.87	81.39	<u>82.30</u> ± 0.72
Mean	73.62	62.35	71.00	68.38	<u>81.49</u>	61.56	72.87	80.62	<b>84.1</b> ± 0.86

Table 11: Anomaly segmentation performance based on Pixel AUROC on MVTecLOCO dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
breakfast_box	<b>89.6</b>	63.61	82.73	87.35	88.53	85.78	76.25	80.76	<u>89.14</u> ± 0.11
juice_bottle	<u>91.37</u>	80.71	86.33	<b>91.99</b>	90.54	90.41	87.06	88.40	89.07 ± 0.10
pushpins	70.66	54.74	<b>82.94</b>	40.72	67.67	41.42	29.42	59.96	<u>75.64</u> ± 0.36
screw_bag	<u>69.94</u>	65.23	58.07	65.35	62.40	67.33	59.74	61.64	<b>71.35</b> ± 0.12
splicing_connectors	63.40	54.16	67.69	<u>71.20</u>	69.71	57.82	56.14	61.02	<b>75.10</b> ± 0.20
Mean	<u>76.99</u>	63.69	75.55	71.32	75.77	68.55	61.72	70.36	<b>80.06</b> ± 0.20

Table 12: Anomaly segmentation performance based on Pixel AUPRO on MVTECLOCO dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
breakfast_box	67.27	36.11	63.8	<b>74.28</b>	<u>73.08</u>	69.67	63.56	58.44	71.36 ± 0.38
juice_bottle	80.75	51.51	77.90	<b>88.78</b>	85.42	84.95	82.88	86.51	<u>87.72</u> ± 0.09
pushpins	61.09	24.68	50.62	52.71	<u>63.52</u>	53.52	59.12	59.25	<b>68.34</b> ± 0.55
screw_bag	54.39	31.27	38.1	61.42	56.12	59.66	<b>71.66</b>	62.45	<u>66.52</u> ± 0.33
splicing_connectors	71.15	56.72	34.77	62.64	67.29	63.62	<u>71.67</u>	68.14	<b>74.70</b> ± 0.25
Mean	66.93	40.06	53.04	67.97	69.09	66.28	<u>69.78</u>	66.96	<b>73.73</b> ± 0.35

Table 13: Anomaly detection performance based on Image AUROC on MPDD dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
tubes	<b>99.64</b>	61.28	89.36	56.48	87.50	89.67	94.47	<u>95.28</u>	93.39 ± 0.57
metal_plate	97.42	80.05	86.92	42.69	<u>99.72</u>	91.87	68.31	<b>100.0</b>	93.81 ± 0.28
connector	<u>99.52</u>	83.33	52.38	86.07	<b>100.0</b>	93.10	<b>100.0</b>	50.00	96.00 ± 0.82
bracket_white	79.89	84.00	50.78	80.33	89.67	83.67	54.55	<u>96.48</u>	<b>100.0</b> ± 0.00
bracket_black	<b>96.48</b>	65.36	62.83	66.69	86.97	50.73	72.63	85.45	<u>93.09</u> ± 0.32
bracket_brown	49.70	70.81	47.89	78.66	<u>95.78</u>	68.70	88.54	85.32	<b>98.08</b> ± 0.14
Mean	87.11	74.14	65.03	68.48	<u>93.27</u>	79.62	79.75	85.42	<b>95.73</b> ± 0.45

Table 14: Anomaly segmentation performance based on Pixel AUROC on MPDD dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
tubes	<b>99.15</b>	76.87	98.44	91.35	98.45	<u>99.08</u>	98.62	98.98	98.55 ± 0.10
metal_plate	<b>98.56</b>	96.23	92.99	91.67	<u>98.30</u>	97.50	93.59	96.52	96.77 ± 0.09
connector	97.38	90.01	92.69	97.93	99.11	98.55	98.68	<u>99.32</u>	<b>99.40</b> ± 0.18
bracket_white	96.74	86.64	90.23	97.21	97.90	<u>98.13</u>	96.63	97.71	<b>98.73</b> ± 0.1
bracket_black	<u>97.68</u>	95.89	94.39	93.79	97.52	96.40	<b>98.42</b>	97.17	97.43 ± 0.21
bracket_brown	95.04	76.11	92.84	95.13	97.15	<u>97.34</u>	<b>98.06</b>	92.46	93.83 ± 2.41
Mean	97.42	86.96	93.60	94.51	<b>98.07</b>	<u>97.83</u>	97.33	97.03	97.45 ± 0.99

Table 15: Anomaly segmentation performance based on Pixel AUPRO on MPDD dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
tubes	<b>96.76</b>	44.84	94.85	71.53	93.83	95.98	95.20	<u>96.27</u>	94.33 ± 0.33
metal_plate	91.53	82.83	74.62	75.47	<b>92.50</b>	<u>92.00</u>	83.99	83.59	90.07 ± 0.22
connector	91.32	72.06	76.98	92.74	96.89	95.29	95.60	<u>97.77</u>	<b>97.98</b> ± 0.60
bracket_white	78.66	69.13	49.65	81.16	83.13	84.71	77.02	<u>93.27</u>	<b>95.33</b> ± 0.37
bracket_black	89.48	93.12	79.65	83.51	<u>93.65</u>	89.14	<b>95.57</b>	89.98	90.17 ± 0.67
bracket_brown	83.62	58.29	85.62	82.69	85.07	<u>94.04</u>	<b>95.37</b>	81.77	84.24 ± 6.38
Mean	88.56	70.04	76.89	81.18	90.84	<u>91.86</u>	90.46	90.44	<b>92.02</b> ± 2.64

Table 16: Anomaly detection performance based on Image AUROC on BTAD dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
01	98.64	80.17	94.46	<u>99.51</u>	98.09	92.23	<u>99.51</u>	94.15	<b>100.0</b> $\pm$ 0.00
02	82.12	65.23	84.27	82.17	81.73	61.73	<u>85.68</u>	75.42	<b>88.5</b> $\pm$ 0.78
03	<u>99.95</u>	74.87	96.3	97.92	<b>100.0</b>	97.65	98.62	95.22	<b>100.0</b> $\pm$ 0.00
Mean	93.57	73.42	91.68	93.20	93.27	83.87	<u>94.60</u>	88.26	<b>96.17</b> $\pm$ 0.45

Table 17: Anomaly segmentation performance based on Pixel AUROC on BTAD dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
01	95.44	59.11	93.05	96.54	95.94	<b>96.98</b>	<u>96.93</u>	64.59	95.86 $\pm$ 0.03
02	94.81	69.29	96.16	95.11	95.18	<b>96.83</b>	<u>96.77</u>	85.67	94.76 $\pm$ 0.88
03	99.55	48.73	99.25	<u>99.56</u>	99.44	<b>99.74</b>	99.12	96.12	99.55 $\pm$ 0.02
Mean	96.60	59.04	96.15	97.07	96.85	<b>97.85</b>	<u>97.62</u>	82.13	96.73 $\pm$ 0.51

Table 18: Anomaly segmentation performance based on AUPRO on BTAD dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
01	66.79	21.57	60.83	75.76	64.34	<u>79.45</u>	<b>83.77</b>	29.75	72.88 $\pm$ 0.12
02	54.32	27.64	<b>67.98</b>	59.19	52.36	<u>66.05</u>	65.58	44.37	55.16 $\pm$ 6.83
03	98.21	18.24	96.99	<u>98.45</u>	97.76	<b>98.92</b>	27.83	88.98	98.18 $\pm$ 0.08
Mean	73.11	22.48	75.27	<u>77.80</u>	71.48	<b>81.47</b>	59.06	54.37	75.41 $\pm$ 3.95

Table 19: Anomaly detection performance based on Image AUROC on VisA dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
candle	<u>94.38</u>	79.43	93.18	86.19	<b>98.59</b>	85.54	89.65	80.52	87.11 $\pm$ 0.29
capsules	69.9	72.77	<u>81.05</u>	61.72	69.92	<b>87.37</b>	76.75	63.73	79.61 $\pm$ 0.72
cashew	94.7	95.5	87.78	90.94	<b>96.29</b>	85.38	93.80	<u>96.11</u>	94.72 $\pm$ 0.16
chewinggum	99.02	83.68	95.18	98.20	<u>99.29</u>	81.92	99.22	98.27	<b>99.49</b> $\pm$ 0.12
fryum	92.98	70.46	92.60	85.06	93.5	77.94	<b>96.58</b>	95.70	<u>95.86</u> $\pm$ 0.14
macaroni1	92.72	72.8	82.48	78.62	91.50	82.06	<u>95.14</u>	<b>95.23</b>	90.66 $\pm$ 0.76
macaroni2	63.44	47.85	69.75	70.05	71.36	81.75	<b>86.25</b>	<u>83.82</u>	82.84 $\pm$ 1.05
pcb1	91.06	72.27	88.07	87.59	<u>95.08</u>	92.60	<b>97.57</b>	93.78	92.92 $\pm$ 0.11
pcb2	79.95	91.17	86.47	83.20	92.46	87.57	91.55	<b>94.95</b>	<u>93.67</u> $\pm$ 0.18
pcb3	82.23	81.29	81.47	72.79	92.46	90.87	<b>97.27</b>	<u>95.92</u>	93.62 $\pm$ 0.16
pcb4	96.29	90.44	95.68	95.67	<u>99.20</u>	96.17	97.62	97.89	<b>99.43</b> $\pm$ 0.03
pipe_fryum	96.54	75.32	96.16	89.28	98.07	85.68	98.36	<u>98.59</u>	<b>99.61</b> $\pm$ 0.11
Mean	87.77	77.75	87.49	83.28	<u>91.48</u>	86.24	85.18	91.21	<b>92.46</b> $\pm$ 0.45



Table 20: Anomaly segmentation performance based on Pixel AUROC on VisA dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
candle	98.75	83.1	97.24	97.68	<u>98.92</u>	<b>99.11</b>	98.41	89.93	97.77 ± 0.03
capsules	96.88	62.39	97.13	90.60	97.62	<b>99.56</b>	<u>99.13</u>	96.93	98.31 ± 0.31
cashew	<u>99.25</u>	74.17	98.57	97.45	98.88	97.23	95.63	98.85	<b>99.49</b> ± 0.02
chewinggum	99.02	84.11	98.83	98.82	98.72	<b>99.37</b>	<u>99.16</u>	98.69	98.10 ± 0.41
fryum	<u>97.08</u>	85.7	93.20	96.20	94.30	96.33	95.45	96.52	<b>97.38</b> ± 0.19
macaroni1	98.71	63.95	98.60	97.85	98.13	99.48	<b>99.73</b>	<u>99.59</u>	99.00 ± 0.13
macaroni2	97.35	79.02	94.65	95.40	96.79	<u>99.33</u>	<b>99.43</b>	98.84	98.20 ± 0.28
pcb1	99.05	27.98	99.29	98.67	99.47	<b>99.65</b>	99.30	98.98	<u>99.61</u> ± 0.01
pcb2	96.40	59.49	97.12	98.12	97.72	<u>98.28</u>	96.13	<b>98.37</b>	98.03 ± 0.09
pcb3	97.23	76.43	97.04	98.06	98.13	<b>98.98</b>	97.99	<u>98.91</u>	98.45 ± 0.05
pcb4	<u>97.97</u>	83.42	97.51	97.00	97.83	<b>98.29</b>	96.58	95.49	95.22 ± 0.27
pipe_fryum	98.79	75.99	98.72	<u>99.19</u>	98.68	98.6	97.97	98.99	<b>99.29</b> ± 0.03
Mean	98.04	71.31	97.32	97.09	97.93	<b>98.68</b>	97.90	97.51	<u>98.24</u> ± 0.20

Table 21: Anomaly segmentation performance based on AUPRO on VisA dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
candle	92.7	80.29	91.65	92.77	94.08	<u>95.30</u>	<b>95.56</b>	77.31	92.49 ± 0.15
capsules	74.64	34.4	81.8	48.42	68.88	<b>92.20</b>	<u>92.09</u>	83.8	82.76 ± 1.18
cashew	<b>93.0</b>	48.33	85.54	82.36	88.01	91.81	90.51	91.57	<u>91.85</u> ± 1.15
chewinggum	<b>89.58</b>	62.66	84.69	84.33	83.86	<u>88.57</u>	85.52	74.87	84.34 ± 1.0
fryum	<u>85.62</u>	71.94	72.39	75.54	78.25	84.8	<b>92.08</b>	82.93	85.47 ± 0.66
macaroni1	89.46	63.37	91.89	88.55	91.74	95.53	<b>97.59</b>	<u>96.06</u>	92.8 ± 0.74
macaroni2	78.74	56.69	71.94	75.76	87.49	<u>94.01</u>	<b>94.23</b>	89.74	88.29 ± 1.89
pcb1	87.24	27.43	85.89	86.39	89.07	<b>95.0</b>	<u>93.55</u>	90.53	90.22 ± 0.28
pcb2	77.83	33.99	77.99	83.68	83.00	<u>89.17</u>	87.26	<b>90.43</b>	84.53 ± 0.53
pcb3	75.03	71.9	71.33	81.37	79.69	90.89	<b>92.48</b>	<u>92.08</u>	85.86 ± 0.38
pcb4	<u>86.53</u>	73.26	83.41	82.47	84.91	<b>89.17</b>	84.38	75.25	73.37 ± 0.81
pipe_fryum	93.14	31.86	81.89	87.99	92.42	<u>94.76</u>	<b>95.46</b>	68.77	93.49 ± 0.08
Mean	85.29	54.68	81.70	80.80	85.12	<b>91.77</b>	<u>91.72</u>	84.45	87.12 ± 0.89

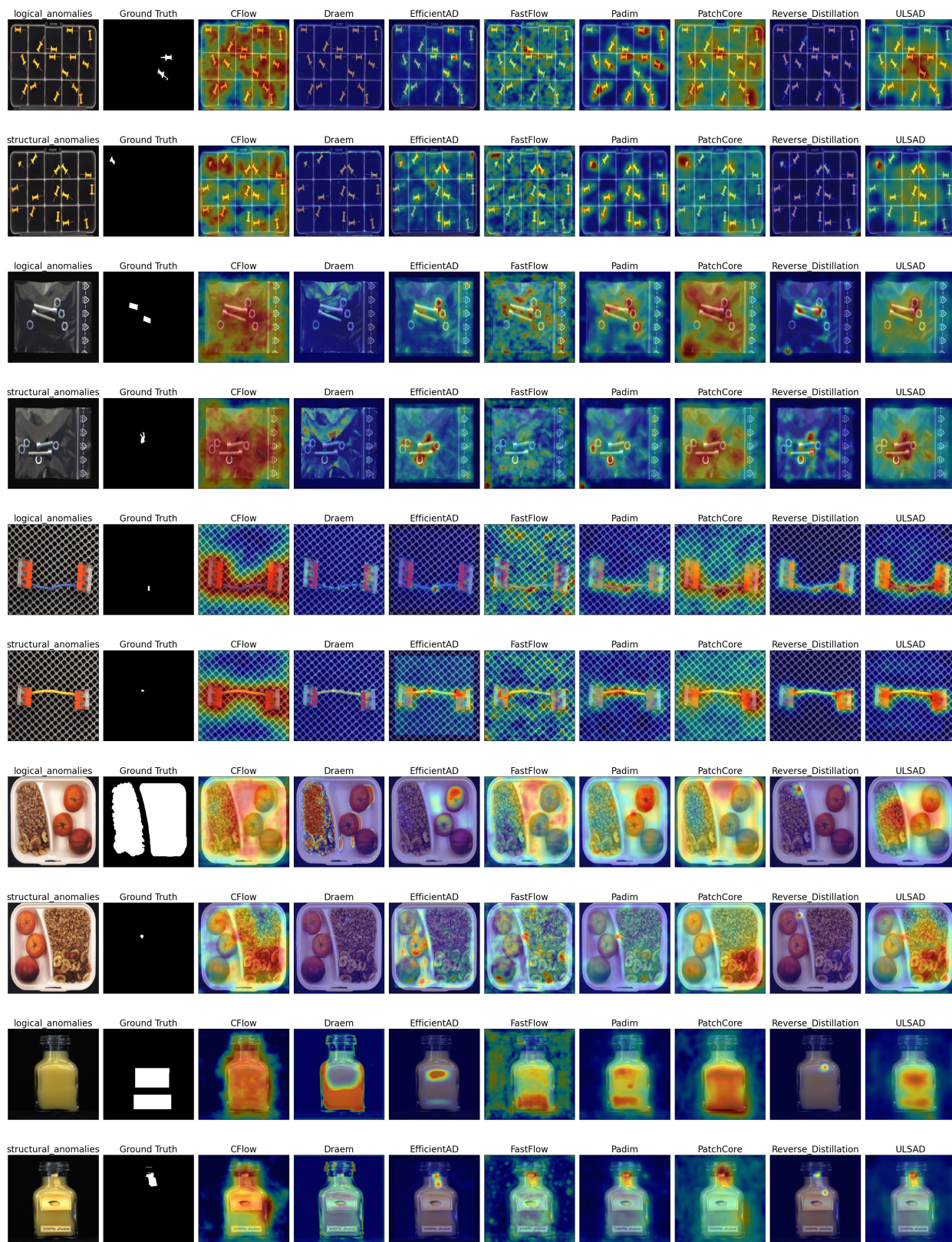


Figure 9: Visualization of anomaly maps on anomalous images from MVTeCLOCO dataset.

### B.1 Performance on MVTecLOCO: Logical and Structural AD

In Tables 22, 23, 24, 25, 26, and 27, we report the anomaly detection and localization results on MVTec LOCO separately for structural and logical anomalies. It can be observed that although PatchCore performs slightly better than ULSAD on structural anomalies, ULSAD delivers competitive results on logical anomalies. Moreover, as discussed previously in Section 5, the improvement in performance with PatchCore comes with three times the memory requirement. Therefore, we consider ULSAD to be an efficient and effective approach for the detection and localization of both logical and structural anomalies.

Table 22: Anomaly detection performance based on Image AUROC on **structural anomalies** of MVTecLOCO dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
breakfast_box	62.3	75.37	71.67	64.17	<b>84.3</b>	48.74	58.67	69.04	<u>81.94</u> ± 0.74
juice_bottle	73.45	52.18	76.52	86.17	<u>96.47</u>	78.29	62.75	<b>99.71</b>	95.53 ± 0.45
pushpins	71.03	66.51	60.17	72.4	77.42	50.17	40.06	<b>92.07</b>	<u>85.17</u> ± 0.65
screw_bag	78.0	69.78	75.07	67.7	<b>86.79</b>	80.35	57.76	82.2	<u>83.01</u> ± 1.56
splicing_connectors	74.35	80.33	70.31	66.85	<u>88.67</u>	63.04	55.84	<b>90.2</b>	80.59 ± 0.35
Mean	71.83	68.83	70.75	71.46	<b>86.73</b>	64.12	55.02	<u>86.64</u>	85.25 ± 0.86

Table 23: Anomaly segmentation performance based on Pixel AUROC on **structural anomalies** of MVTecLOCO dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
breakfast_box	<u>94.72</u>	61.6	84.86	88.41	94.31	91.75	82.79	67.36	<b>95.35</b> ± 0.16
juice_bottle	88.81	83.75	85.67	90.73	<u>95.98</u>	89.57	80.74	<b>97.24</b>	89.67 ± 0.23
pushpins	<u>91.4</u>	46.14	85.54	90.82	<b>95.16</b>	87.83	48.52	89.48	88.38 ± 0.22
screw_bag	95.37	72.27	91.74	94.97	97.45	<b>97.78</b>	66.67	97.37	<u>97.64</u> ± 0.24
splicing_connectors	98.3	93.03	95.79	96.25	<b>98.82</b>	97.25	91.82	<u>98.55</u>	98.3 ± 0.1
Mean	93.72	71.36	88.72	92.24	<b>96.34</b>	92.84	74.12	90.0	<u>93.87</u> ± 0.2

Table 24: Anomaly segmentation performance based on Pixel AUPRO on **structural anomalies** of MVTecLOCO dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
breakfast_box	70.0	41.92	71.29	<b>81.97</b>	<u>78.71</u>	76.67	26.67	65.87	75.43 ± 0.67
juice_bottle	80.56	50.44	87.61	<u>92.64</u>	<b>95.1</b>	90.81	61.67	92.08	90.1 ± 0.31
pushpins	68.34	20.04	61.72	70.01	<u>75.88</u>	62.63	45.30	<b>78.62</b>	68.34 ± 0.84
screw_bag	82.07	39.41	73.69	85.1	88.92	<b>93.77</b>	53.95	91.38	<u>91.98</u> ± 0.73
splicing_connectors	87.75	68.61	67.8	82.35	<u>91.69</u>	83.57	63.44	<b>94.19</b>	90.84 ± 0.35
Mean	77.74	44.08	72.42	82.41	<b>86.06</b>	81.49	50.21	<u>84.43</u>	83.34 ± 0.62

Table 25: Anomaly detection performance based on Image AUROC on **logical anomalies** of MVTecLOCO dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
breakfast_box	76.86	68.71	75.81	62.0	<u>83.43</u>	55.27	61.55	<b>83.54</b>	83.33 ± 1.63
juice_bottle	80.52	70.09	80.89	92.94	94.61	76.4	77.07	<b>99.12</b>	<u>98.82</u> ± 0.14
pushpins	71.15	37.9	58.81	50.72	<u>74.53</u>	52.98	60.27	72.0	<b>85.23</b> ± 0.72
screw_bag	60.72	52.34	<u>64.01</u>	55.29	59.36	55.09	63.64	58.8	<b>66.33</b> ± 0.9
splicing_connectors	67.5	56.15	74.54	69.43	<u>80.27</u>	61.74	53.10	75.64	<b>86.27</b> ± 1.33
Mean	71.35	57.04	70.81	66.08	<u>78.44</u>	60.3	63.13	77.82	<b>84.0</b> ± 1.07

Table 26: Anomaly segmentation performance based on Pixel AUROC on **logical anomalies** of MVTecLOCO dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
breakfast_box	<u>90.43</u>	65.5	85.35	89.23	90.32	87.35	74.55	85.58	<b>91.41</b> ± 0.1
juice_bottle	92.56	80.49	90.23	<b>95.29</b>	<u>93.97</u>	92.97	86.95	91.98	93.74 ± 0.08
pushpins	70.38	56.06	<b>82.91</b>	41.95	69.39	41.86	67.93	61.12	<u>75.96</u> ± 1.11
screw_bag	67.75	65.43	58.93	66.02	63.81	68.4	<b>75.75</b>	61.85	<u>72.3</u> ± 0.41
splicing_connectors	60.8	52.1	66.82	<u>69.83</u>	68.06	55.5	57.67	59.01	<b>74.19</b> ± 0.19
Mean	76.38	63.92	76.85	72.46	<u>77.11</u>	69.22	72.57	71.91	<b>81.52</b> ± 0.54

Table 27: Anomaly segmentation performance based on Pixel AUPRO on **logical anomalies** of MVTecLOCO dataset.

Category	CFLOW (2021)	DRÆM (2021)	FastFlow (2021)	PaDiM (2021)	PatchCore (2022)	RD (2022)	DFR (2020)	EffAD (2024)	ULSAD (Ours)
breakfast_box	68.79	32.23	64.39	69.74	<u>72.27</u>	68.77	43.60	52.19	<b>73.97</b> ± 0.8
juice_bottle	79.67	52.72	77.9	<u>91.69</u>	87.88	84.22	62.87	87.82	<b>91.36</b> ± 0.12
pushpins	59.07	26.21	47.71	51.93	<u>63.47</u>	53.84	41.80	58.36	<b>68.78</b> ± 0.98
screw_bag	50.7	26.98	25.58	<u>53.92</u>	46.8	48.72	54.59	52.8	<b>61.48</b> ± 0.45
splicing_connectors	<u>65.86</u>	53.71	26.66	57.66	62.21	58.85	34.72	62.7	<b>72.66</b> ± 0.27
Mean	64.82	38.37	48.45	64.99	<u>66.53</u>	62.88	47.52	62.77	<b>73.65</b> ± 0.61

## C Extended Ablations

In this section, we provide additional ablations on the local branch in Table 28 and the total architecture in Table 29. Lastly, in Table 30 we provide the per-category results for the ablation on the pre-trained backbone which is summarized in Figure 7.

Table 28: Ablations for local branch. Style: I-AUROC | P-AUROC | P-AUPRO.

category	$\lambda_l = 0$	$\lambda_l = 0.01$	$\lambda_l = 0.5$	$\lambda_l = 0.9$	$\lambda_l = 1.0$
breakfast_box	78.64   <u>88.28</u>   <u>74.22</u>	<u>79.2</u>   <b>88.51</b>   <b>74.35</b>	77.86   87.79   71.27	77.95   86.89   67.06	<b>79.44</b>   86.96   65.36
juice_bottle	<b>97.82</b>   <u>92.14</u>   89.24	<u>97.76</u>   <b>92.23</b>   89.38	97.56   88.78   88.16	97.36   84.39   84.63	97.08   83.61   83.47
pushpins	72.4   69.81   <u>65.69</u>	72.77   69.84   65.68	<b>79.92</b>   <b>74.49</b>   <b>69.03</b>	<u>76.98</u>   <u>74.35</u>   63.17	76.53   73.69   65.18
screw_bag	66.42   66.6   64.39	67.18   68.47   <u>65.92</u>	<b>68.06</b>   69.13   <b>66.22</b>	<u>67.56</u>   <b>69.33</b>   63.61	66.34   <u>69.31</u>   62.09
splicing_connectors	<b>73.05</b>   59.04   73.3	<u>72.84</u>   59.15   73.29	72.29   62.66   72.39	72.79   <u>64.33</u>   70.74	72.36   <b>64.5</b>   70.57
Mean	77.67   75.17   73.37	77.95   75.64   <b>73.72</b>	<b>79.14</b>   <b>76.57</b>   <u>73.41</u>	78.53   75.86   69.84	78.35   75.61   69.33

Table 29: Ablations for total architecture. Style: I-AUROC | P-AUROC | P-AUPRO.

category	$\mathcal{L}_{pg}^d$	$\mathcal{L}_{pg}^d; \mathcal{L}_{lg}$	$\mathcal{L}_{pg}$	$\mathcal{L}_{pg}; \mathcal{L}_{lg}$
$\lambda_l = \lambda_g = 0.0$				
breakfast_box	77.29   <b>90.41</b>   77.16	82.82   89.85   76.92	66.01   87.36   67.5	<b>85.08</b>   <u>90.19</u>   75.36
juice_bottle	96.48   <b>92.01</b>   88.82	<b>97.93</b>   91.82   <b>89.26</b>	91.2   91.98   85.38	<u>97.29</u>   <u>92.0</u>   <u>89.18</u>
pushpins	70.89   80.89   70.67	<b>78.66</b>   <b>88.09</b>   <b>79.11</b>	<u>74.67</u>   77.37   58.75	74.61   <u>85.86</u>   <u>76.33</u>
screw_bag	<u>65.48</u>   65.87   64.04	63.02   <u>68.13</u>   <b>65.51</b>	61.14   59.51   57.75	<b>66.93</b>   <b>68.67</b>   <u>65.35</u>
splicing_connectors	78.29   69.68   75.61	<b>84.55</b>   <u>72.71</u>   <b>76.54</b>	65.31   53.4   66.74	<u>81.5</u>   <b>73.11</b>   <u>76.01</u>
Mean	77.69   79.77   75.26	<b>81.4</b>   <b>82.12</b>   <b>77.47</b>	71.67   73.92   67.22	<u>81.08</u>   <u>81.97</u>   <u>76.45</u>
$\lambda_l = \lambda_g = 0.5$				
breakfast_box	79.22   <b>90.87</b>   <b>78.32</b>	82.45   88.49   71.03	71.36   87.64   67.38	<b>83.36</b>   <u>89.34</u>   <u>72.36</u>
juice_bottle	96.3   <u>91.17</u>   <b>88.89</b>	<b>98.08</b>   87.06   <u>88.05</u>	91.14   <b>91.84</b>   85.92	<u>97.46</u>   88.81   87.71
pushpins	79.86   <u>84.89</u>   <u>77.97</u>	<u>82.46</u>   <b>87.62</b>   <b>80.49</b>	78.64   81.14   65.12	<b>88.07</b>   74.22   66.45
screw_bag	<u>66.58</u>   68.83   <u>66.11</u>	65.11   <u>70.04</u>   62.81	62.09   67.32   62.92	<b>70.6</b>   <b>71.58</b>   <b>67.01</b>
splicing_connectors	80.53   <u>73.47</u>   <b>75.44</b>	<b>82.85</b>   73.03   <u>75.13</u>	69.33   55.02   63.69	81.27   <b>74.94</b>   74.89
Mean	80.5   <b>81.85</b>   <b>77.35</b>	<u>82.19</u>   <u>81.25</u>   <u>75.5</u>	74.51   76.59   69.01	<b>84.15</b>   79.78   73.68

Table 30: Ablations for backbone on MvTec-LOCO. Style: I-AUROC | P-AUROC | P-AUPRO.

Class	ResNet50	ResNet152	Wide-ResNet50-2	Wide-ResNet100-2
breakfast_box	82.41   <u>89.74</u>   <u>73.09</u>	<b>85.11</b>   <b>91.15</b>   72.0	<u>84.46</u>   89.18   72.21	82.37   89.25   <b>74.02</b>
juice_bottle	96.9   <u>92.23</u>   89.38	<u>97.64</u>   91.66   <u>89.78</u>	97.11   88.9   87.94	<b>98.74</b>   <b>92.87</b>   <b>91.07</b>
pushpins	81.79   <b>79.49</b>   <b>75.15</b>	73.28   76.49   63.71	<b>85.48</b>   75.46   67.82	<u>82.48</u>   <u>76.55</u>   <u>70.67</u>
screw_bag	66.46   <u>69.14</u>   <b>67.01</b>	68.06   68.63   <u>66.97</u>	<u>71.14</u>   <b>71.57</b>   66.82	<b>73.1</b>   68.99   66.88
splicing_connectors	80.88   <u>72.76</u>   74.39	<u>83.53</u>   <u>76.31</u>   77.24	82.59   75.21   75.05	<b>84.71</b>   <b>79.95</b>   <b>78.65</b>
Mean	81.79   79.49   <u>75.15</u>	81.52   <u>80.85</u>   73.94	<u>84.16</u>   80.06   73.97	<b>84.28</b>   <b>81.52</b>   <b>76.26</b>