



A Generative Encoder-Decoder Model for Automated Quality Control Inspections

Khedidja Mekhilef¹, Fayçal Abbas² and Mounir Hemam¹

¹University Abbas Laghrour- Khenchela, ICOSI Laboratory, BP 1252 El Houria, 40004, Algeria

²University Abbas Laghrour- Khenchela, LESIA Laboratory, BP 1252 El Houria, 40004, Algeria

Received 16 April 2024, Revised 13 November 2024, Accepted 28 November 2024

Abstract: This paper introduces a novel generative model based on an encoder-decoder architecture for defect detection within Industry 4.0 frameworks, focusing on the escalating need for automated quality control in manufacturing settings. Precision and efficiency, crucial in such environments, are significantly enhanced by our approach. At the core of our methodology is the strategic incorporation of random Gaussian noise early in the image processing sequence. This deliberate interference disrupts the model's ability to reconstruct images of defective parts, thereby enhancing both the accuracy and robustness of defect detection.

The model further integrates skip connections during the decoding phase, with a special emphasis on the first two connections. These are augmented with multi-head attention mechanisms and spatial reduction techniques, followed by targeted convolutions. This intricate configuration helps preserve vital local features while filtering out superfluous data, facilitating precise image reconstruction and effectively addressing the often problematic issue of locality loss during the upsampling process. Moreover, our model excels in maintaining contextual integrity and capturing multi-scale features, which is crucial for detailed defect detection. Each block of the architecture connects to a scaled version of the original image, allowing for nuanced feature analysis. Extensive testing and validation on real-world datasets have proven the model's high efficiency and accuracy in identifying defects, marking a significant advancement in automated quality control systems.

Keywords: Anomaly Detection, Vision Transformer, Quality Control, Industry 4.0

1. INTRODUCTION

Advanced technologies such as robotics, artificial intelligence, machine vision, big data, cloud computing, and machine learning have revolutionized manufacturing. They have given rise to what is known as Industry 4.0, which has played a major role in the application of automated visual inspection. This has helped avoid many problems that can be caused by human inspection by using artificial vision techniques, such as cameras and capture devices, to record images and transfer them to a machine to check product quality [1][2][3][4].

The quality of industrial products is defined by their compliance with established standards. Any defect impacting product quality indicates it has not met the required standards, leading to potential issues such as safety risks, breakdowns, material damage, or even injuries. These incidents can result in financial losses for companies and a negative reputation. This is why defect detection is fundamental in product quality control. Defect detection is the process of identifying anomalies that occur during

production, such as contamination, scratches, cracks, color changes, etc. Computer vision is one of the most widely adopted fields for this task. It involves capturing images of the product, with and without defects, and then letting the model operate until it can distinguish between the two. This produces meticulous results. Since defects can vary in different ways, annotating all types of defects becomes impossible due to the time required.

This has prompted researchers to focus on unsupervised learning. Some have explored methods based on feature integration. The fundamental idea of this approach is to generate, during training, a significant vector space to represent normal data. During the testing phase, results are compared to this vector to classify whether they indicate a defect or not. Conversely, other researchers have opted for reconstruction approaches. The main idea behind reconstruction is to train the model exclusively on images without defects. Although this approach creates divergences when processing images containing defects, these differences effectively reveal the presence of defects



during the testing phase. Most work has adopted the reconstruction approach based on convolutional layers, incorporating architectures such as Autoencoder Networks [5][6][7], GAN (Generative Adversarial Networks) [8][9][10], and Variational Autoencoders [11][12][13]. However, the major drawback of these convolutional layers lies in their excessive focus on locality, which limits their explicit modeling of long-term dependencies. This limitation results in often imperfect reconstruction, even for non-defective images during the testing phase, thus compromising accurate defect detection.

With the emergence of the Vision Transformer architecture [14], inspired by the Natural Language Processing (NLP) model [15] known for efficiently modeling long-term dependencies, several research studies have been encouraged to adopt this architecture in reconstruction-based methods for defect detection. Some have even substituted the autoencoder's encoder with a transformer [16][17], while others have explored using it to create a self-attention based autoencoder for feature reconstruction [18]. However, despite successful modeling of the global context by this architecture, its use has sometimes led to a lack of locality [19]. Mathian et al. [20] aimed to combine locality and globality by using an autoencoder composed of a sequence of a convolutional layer followed by a self-attention mechanism. However, this raises concerns about the quality of the extracted locality, as the exclusive use of convolutional layers may require a well-defined sequence for efficient extraction. Considering the inherent visual complexity of images, characterized by intricate patterns and details, accurate reconstruction of images or features requires considering both global and local information. Nevertheless, a challenge persists in the context of reconstruction-based methods during the testing phase, where the presence of defective images can lead to the reconstruction of defects, thereby complicating the precise detection and localization of anomalies.

In this paper, to fully leverage the complementarity of local and global features, an encoder-decoder architecture is proposed. The initial layers of the encoder capture texture features, while the final layers focus more on semantic features. To restore the image from the extracted features, the decoder applies a set of upsampling and convolution operations. However, during the decoding phase where upsampling operations are performed, there can be a loss of locality. To enable the decoder to fully utilize this information for precise image reconstruction, inspired by U-Net [21], it is integrated with skip connections, the first two of which are combined with multi-head attention, followed by spatial reduction inspired by [22], and then convolution aimed at retaining only specific local features and eliminating those that are not necessary. To maintain the integrity of contextual information on one hand and capture features at different spatial scales on the other, each block is associated with an equivalent representation of the original image, but at a reduced scale. To prevent the problem of reconstructing the defect and hinder the reconstruction of

the defective part from random Gaussian noise, the latter is added at the beginning of the image. In addition to this, to enrich the dataset dedicated to defect detection and localization, a new class of data is created. The remainder of this paper is structured as follows: Section 2 provides an overview of the related work, Section 3 describes the proposed method, Section 4 presents the experiments, section 5 Describe integration in real time 6 present the limitations of the model. and Section 7 concludes the paper.

2. RELATED WORK

A. Methods Based on Reconstruction

Since the database contains only non-defective images, some research has explored the effectiveness of CNNs in reconstruction methods. Bergmann et al. [5] introduced structural similarity as a metric, replacing the simple pixel difference (L2) in their approach. Yang et al. [23] introduced the concept of multi-sequence by combining model blocks at different scales. Zavrtnik et al. [24] proposed image inpainting, masking specific portions in the images to prompt the model to reconstruct the defective parts as if they were non-defective. Zhou et al. [25] Based their approach on the difference between the structural information of the original image and the reconstructed image to detect defects. Li et al. [26] introduced the concept of superpixels to divide the image into regions, and then masked these regions randomly to prevent the reconstruction of defects in the test portion. Hou et al. [6] introduced the concept of multi-scale block-wise memory in autoencoders to maximize the difference between the reconstruction of defective and non-defective images. Jiang et al. [27] introduced an 'Interpretability-Aware' loss in the autoencoder to enhance result interpretability during training and testing. Li et al. [28] introduced continual learning. Zhao et al. [29] introduced efficient channel attention, as well as a strategy to better distinguish the foreground from the background. Li et al. [30] used the Dual Attention mechanism to optimize reconstruction, with a loss function aimed at enhancing defect detection. Zhang et al. [31] proposed AGUR-Net with EfficientNet-B2 as the encoder and an Atrous Spatial Pyramid Pooling module in the decoder, integrating a residual fusion with attention and a dual-threshold segmentation method to enhance defect detection. Chen et al. [32] proposed a 3D model with three components: a multimodal reconstruction to restore the normal image, segmentation to extract defects, and an attention model to enhance anomaly detection. Wang et al. [33] analyzed the image in small patches to better locate anomalies.

Other works have explored the potential use of transformers to enhance data representation. Lee et al. [17] introduced the transformer as an encoder for the CNN autoencoder. De et al. [34] applied masking to hide information, focusing particularly on the masking of patches inside blocks. You et al. [35] introduced the transformer into a method based on feature reconstruction. Mishra et al. [16] introduced a Gaussian mixture density network to model the distribution of representative vectors generated

by the encoder of the Vision Transformer in the context of defect detection and localization. Yang et al. [36] developed an autoencoder with a Vision Transformer encoder to capture global information and a memory model to store normal data. A Coordinate Attention block enhances the representation before the decoder reconstructs the final image. Shang et al. [37] replaced the autoencoder's encoder with a transformer featuring a "defect-aware" mechanism and graph-based positional encoding to enhance performance. Our model is based on the overarching concept of reconstruction. It is inspired by previous work that has demonstrated the effectiveness of CNNs for extracting local features, as well as transformers for capturing global features. We have also evaluated multi-scale approaches, which are recognized for their relevance in capturing details at different levels of granularity.

B. Feature Integration based Methods

To improve the performance of unsupervised methods, some approaches strive to incorporate the idea of a representative vector or vector space.[38][39][40][41] Create a hypersphere space during training by minimizing the distance between normal points and the center of the hypersphere. During the test phase, if the distance is no longer close or identical, the instance is considered defective.[42][43] opt for the use of normalized vectors through distribution estimation methods. During the test stage, if the distance between the normal and observed distribution is higher, the instance is considered defective.[44][45] follows an approach where the teacher is considered a reference vector. During training, a student model tries to adapt to this teacher. During the test phase, if the student fails to mimic the teacher, the instance is considered defective. Unlike previous work, our model does not rely on their main idea. Nevertheless, some of these related studies have applied fine-tuning, and we have leveraged their results to integrate a pre-trained model.

3. METHOD

A. Feature Extraction

Pre-trained CNNs are recognized as being among the most effective models for producing discriminative features that have a significant impact on tasks of defect detection and localization [23], as shown in Figure 3, To enable the model to capture local information, the first three blocks of the VGG19 Network pre-trained on the ImageNet dataset are used, the first and second are designed for the extraction of texture features, while the third acts as an intermediary between texture and semantic information.

$f \in \mathbb{R}^{H \times W \times C}$ The feature map of the last block where C and $H \times W$ indicate the channel and the spatial dimensions of the feature map, respectively. As the vision transformer set of tokens unfolds, the feature map is initially divided into a set of tokens $N = \frac{H \times W}{p^2}$ where $P \times P$ represents the resolution of each token, these tokens are linearly projected into latent vectors of size D combined with position encoding to restore the information to its position before they are introduced to the transformer block to model global information as the

permutation is invariant. Regarding The transformer block, it follows the structure of the classic architecture shown 1, the encoding passes through a multi-head attention mechanism in the first sub-block and a forward propagation in the second sub-block, and normalization and residual connection in both sub-blocks. Everything related to feature extraction is illustrated in Figure 2.

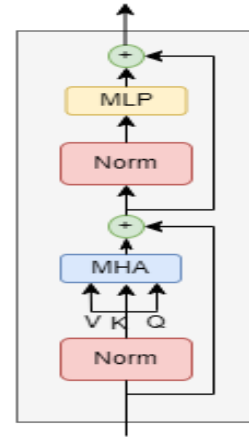


Figure 1. Block transformer components.

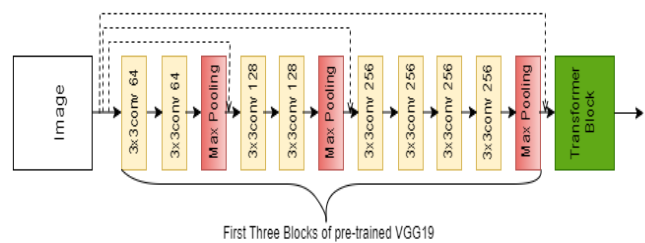


Figure 2. The encoder architecture.

B. Image Reconstruction

In order to reconstruct the image and decipher hidden features, we begin by reshaping the dimensions of the output of the hybrid encoder, changing from $N \times D$ to $n_1 \times n_2 \times D$, where $n_1 = \frac{H}{P}$ and $n_2 = \frac{W}{P}$. This is followed by the application of a convolutional layer to restore the original dimensions, and then a series of upsampling operations with a magnification factor of 2 to enhance spatial resolution, and conv3x3 for extracting more complex features. A Relu activation is also applied. In addition, a sigmoid function is used on the final results to normalize the values between 0 and 1. The structure of this proposed model, designed according to an encoder-decoder scheme, naturally allows for the integration of 'skip' type connections between the encoder and decoder. These connections are crucial for effectively associating high-resolution local information with low-resolution global information.

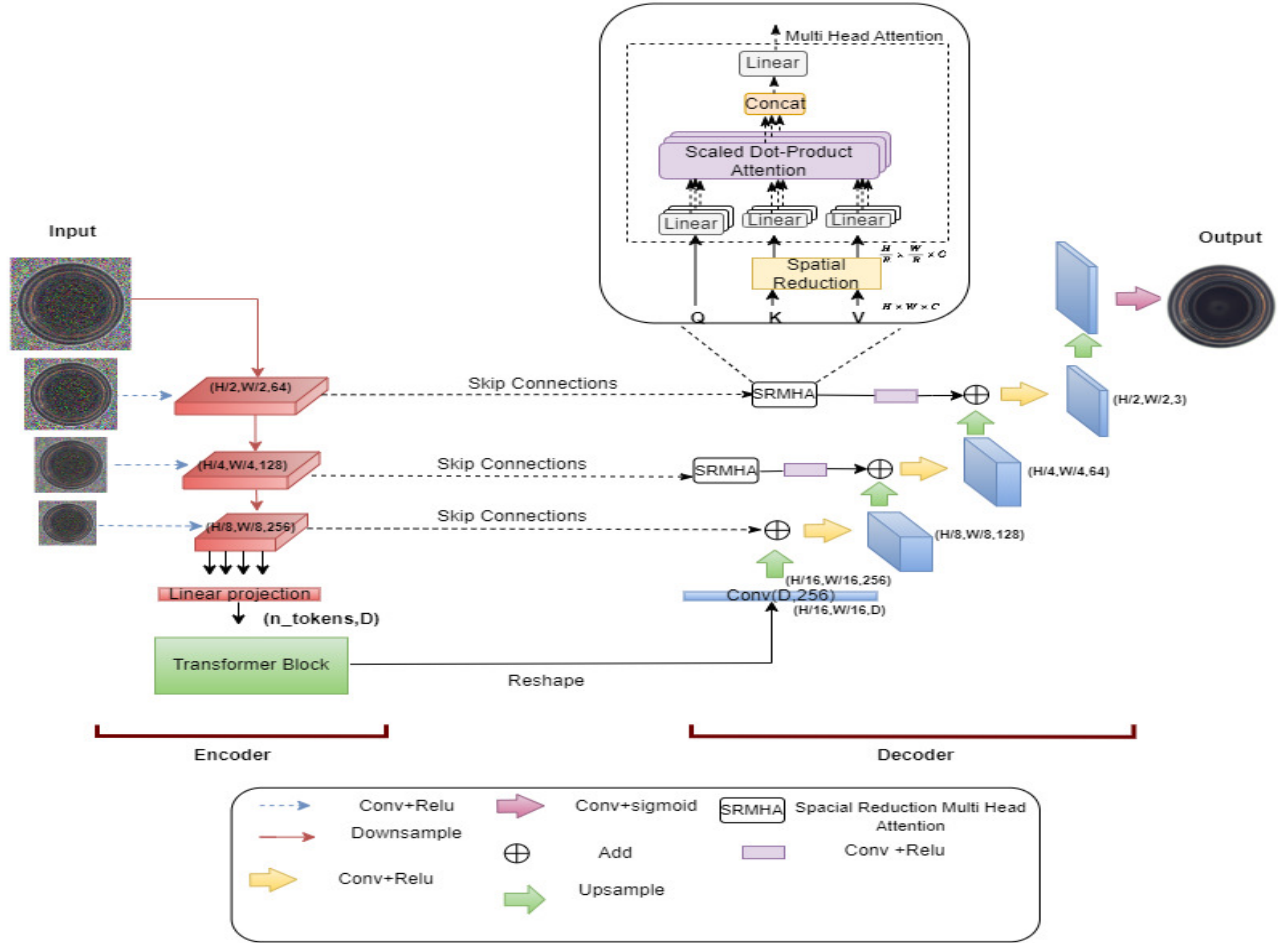


Figure 3. The overall architecture of our model.

As the encoder becomes more complex, the processed information becomes more global and elaborate, which can, however, lead to a loss of local details during the decoding process. This particularly affects the reconstruction of objects with variable structures and complex patterns. To address this problem, multi-head attention combined with spatial reduction is integrated at the level of the first and second residual connections. This approach allows for a significant improvement in the ability to weigh and integrate texture information across the entire image, thereby enhancing the representation of relevant features in the overall context of the scene. This is succeeded by a conv3x3 and a Relu activation function to accentuate local details.

The multi-head attention mechanism (MHA), designed to identify distant interdependencies, operates as follows: the linear projections of keys (K), queries (Q), and values (V), all of which have the same dimension, are distributed across multiple heads. In each head, a multiplication is per-

formed between the keys and queries, after which a softmax function is applied to the result of this multiplication. The resulting output is then adjusted by multiplying it with the corresponding values. This process can be expressed in the following way.

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (1)$$

The first and second blocks create a high-resolution feature map, whose integration into MHA increases the computational load and memory usage. The implementation of Spatial Reduction for multi-head attention involves adjusting the dimensions of the keys and values via a spatial reduction R, before proceeding to the attention operation 1.

$$K_{reduced} = \text{Reshape}\left(\frac{HW}{R^2}, C \cdot R^2\right)W(C \cdot R^2, C) \quad (2)$$

$$V_{reduced} = K_{reduced} = \text{Norm}(K_{reduced}) \quad (3)$$

W refers to a linear projection designed to preserve channel dimensions, while Norm refers to the normalization layer. To enrich the information represented by the three blocks of the pre-trained CNN and to address issues related to resolution reduction and capturing features at various spatial scales, additional features are introduced for each block. These features are generated from a sequence that includes a 2×2 average pooling operation, a 2×2 convolution, followed by activation using the Relu function on the original image. Careful restoration of visual data can sometimes lead to the reappearance of defects during the testing phase. To avoid this, a proactive approach has been implemented: the deliberate introduction of a random Gaussian disturbance in the input image. This disturbance is designed to mask certain information while preserving the overall quality of the reconstruction. The formula used for integrating this noise is as follows:

$$X_{noisy} = X + \lambda \quad \text{where, } \lambda \sim N(0, \sigma^2) \quad (4)$$

As X represents the input image and σ is the maximum standard deviation of the Gaussian noise added to the input image.

C. The Loss Function

During the training phase, a loss function was used that combines both the pixel-focused L_2 method and the SSIM [5]. The pixel-focused L_2 method assesses the error in the value of each corresponding pixel, while the SSIM evaluates brightness, defined as the average value of all the pixels; contrast, measured by the standard deviation of the pixel intensities; and structural similarity, which indicates the correlation between the two images and is measured by the divergence in intensity distributions.

$$\text{SSIM}(X, \hat{X}) = \frac{(2\mu_x\mu_{\hat{x}} + C_1)(2\sigma_x\sigma_{\hat{x}} + C_2)}{(\mu_x^2 + \mu_{\hat{x}}^2 + C_1)(\sigma_x^2 + \sigma_{\hat{x}}^2 + C_2)} \quad (5)$$

$$\text{Loss}_T = L_2(X, \hat{X}) + \text{SSIM}(X, \hat{X}) \quad (6)$$

where,

- X the original image.
- \hat{X} the reconstructed image.
- μ_x the average sample of the image X .
- $\mu_{\hat{x}}$ the average sample of the image \hat{X} .
- σ_x^2 the variance of X .
- $\sigma_{\hat{x}}^2$ the variance of \hat{X} .
- $\sigma_x\sigma_{\hat{x}}$ the covariance of X and \hat{X} .
- $C_1 = (K_1L)^2$ and $C_2 = (K_2L)^2$ are two variables to stabilize the division with a weak denominator.
- L is the dynamic range of pixel values (typically it's $2^{\text{bits per pixel}} - 1$).
- $K_1=0.01$ and $K_2 = 0.03$ by default.

During the test phase and to evaluate the performance of our model, we use the multi-scale gradient magnitude similarity method (MSGMS) [24], a multi-scale extension

of GMS [46], to evaluate the similarity of structure and contours, in conjunction with the L_2 loss to calculate the pixel-to-pixel difference between the values. This method allows us to estimate the anomaly score between the reconstructed image and the original image. The function for calculating this anomaly score is presented as follows:

$$A_{\text{score}} = (1_{H \times W} - \text{MSGMS}(X, \hat{X}))\text{Conv}_f + L_2(X, \hat{X})\text{Conv}_f \quad (7)$$

The anomaly score is obtained by subtracting the anomaly map obtained from MSGMS from $1_{H \times W}$, where $1_{H \times W}$ is a matrix of ones, and then adding the result to the anomaly map obtained from the L_2 loss. The anomaly maps obtained from MSGMS and L_2 have been previously processed by a mean filter convolution of size 21 Conv_f .

A_{score} is a matrix representing the anomaly score of each pixel. To calculate the score for the entire image, the maximum among all scores is taken into account.

4. EXPERIMENTS

A. Data Sets

In the context of our study, we used two datasets to assess the effectiveness and accuracy of our model for defect detection.

The first dataset, as illustrated in Figure 4, consists of real images of buttons that we created. This set contains 173 images, divided into two categories: 131 images for training and 42 for testing, with 14 non-defective and 28 defective images. Each image in this dataset has dimensions of 704 pixels in width by 708 pixels in height and is presented in RGB color format. These images were captured using a mobile phone camera, which offers a high resolution of 4032x2268 pixels. This capture method guarantees high image quality, essential for detailed and precise analysis. To enhance the effectiveness of the defect detection process, masks were generated for all images showing anomalies. These masks play a crucial role in our study, as they allow precise localization of defects on the images. This method greatly facilitates the evaluation of our model's performance in terms of detection and localization of defects on button images.

The second dataset is MVTec Anomaly Detection (MVTec AD) [47], a diverse and specialized dataset crucial for assessing how well anomaly detection techniques work in unsupervised machine learning. This dataset consists of fifteen distinct industrial categories, including five different texture types and ten different object categories. All of the categories are in RGB format, except for the "grid," "zipper," and "screw" categories, which are in grayscale format. This variety allows for a comprehensive and exhaustive evaluation of anomaly detection models, offering a wide range of possible scenarios and use cases. Detailed information on the dataset including the distribution of categories is provided in Table I. For each category, MVTec AD provides two distinct sets of images: one is used for testing, and another for training. These training pictures are carefully selected to present no defects, ensuring that the

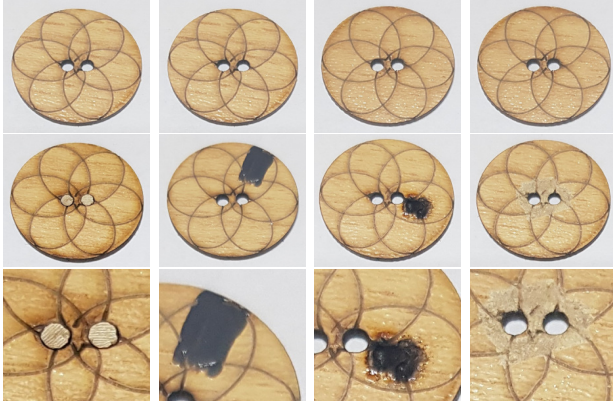


Figure 4. Non-defective samples (first row) and Defective samples (second row) with Defect overview (last row).

models learn from examples without anomalies. In contrast, the test set consists of both defective and non-defective images. This composition is crucial for testing the models' ability to distinguish anomalies from normal conditions in real environments. A particularly relevant aspect of MVTEC AD is the inclusion of annotated masks for each defective image in the test set. These masks provide precise information on the location and nature of defects in the images. Using these masks not only allows for assessing whether a model can detect an anomaly but also measures its accuracy in locating and characterizing specific defects.

TABLE I. Detailed Information on the MVTEC AD Dataset: Learning Set (L. Set), Evaluation Set (ES), Real (R) Defective (D), Defect Types (dfts).

Class	L. Set	ES(R,D)	dfts	Res
Bottle	209	(20,63)	3	900 × 900
Cable	224	(58,92)	8	1024 × 1024
Capsule	219	(23,109)	5	1000 × 1000
Hazelnut	391	(40,70)	4	1024 × 1024
Metal nut	220	(22,93)	4	700 × 700
Pill	267	(26,141)	7	800 × 800
Screw	320	(41,119)	5	1024 × 1024
Toothbrush	60	(12,30)	1	1024 × 1024
Transistor	213	(60,40)	4	1000 × 1000
Zipper	240	(32,119)	7	1024 × 1024
Carpets	280	(28,89)	5	1024 × 1024
Grid	264	(21,57)	5	1024 × 1024
Leather	245	(32,92)	5	1024 × 1024
Tile	230	(33,84)	5	840 × 840
Wood	247	(19,60)	5	1024 × 1024

B. Implementation Details

At the beginning of the process, before feature extraction begins, images are first scaled to 224 pixels. Then, the parameters of the transformer block head and the multi-head attention for the second and first levels of the skip

connection are set to 4, 2, and 1, respectively. Furthermore, the encoding parameters of the transformer block size (D) and the multi-head attention for the second and first skip connections are fixed at 512, 128, and 64. Lastly, the spatial reduction rate for the first skip connection is set to 8, and for the second skip connection, it is fixed at 4. Dropout with a value of 0.25 is applied in both the MLP and the attention blocks of the transformer. The model is run with the Adam optimizer with a learning rate equal to 0.0001. The dataset is divided into 80% for processing and 20% for validation with batch sizes of 8. The model was trained for 2000 epochs, with an early stopping mechanism activated from epoch 800. This mechanism ends the training if the validation loss shows no improvement for 300 consecutive epochs. It is important to note that the training loss and validation loss remain very close in value, with minimal difference. This suggests that the model generalizes well to the validation data and shows no significant signs of overfitting. Regarding the noise rate, each class is run and evaluated individually and independently of other categories. We ran each class with a different noise rate to select the optimal rate offering the best performance. The noise rates chosen for each class are as follows: 0.1 for 'bottle', 0.25 for 'cable', 0.2 for 'capsule', 0.2 for 'hazelnut', 0.3 for 'metal nut', 0.09 for 'pill', 0.2 for 'screw', 0.3 for 'toothbrush', 0.3 for 'transistor', 0.1 for 'zipper', 0.4 for 'carpet', 0.16 for 'grid', 0.1 for 'leather', 0.11 for 'tile', 0.09 for 'wood', and 0.09 for the new class 'constructed button'. The network was implemented in PyTorch with GPU RTX 3050 6G.

C. Results and Discussion

In order to assess how well our model finds and detects flaws, we undertook an extensive comparison of our results with those obtained by current state-of-the-art methods in this field. Our analysis focused on various approaches, including knowledge distillation with KDAD [44], various reconstruction methods such as SMAI [26], AnoGAN [8], VTBA [36], GAP [33], AESSIM [5] and HaloAE [20], as well as one class classification techniques like FCDD [40], specifically for defect localization. Additionally, we also examined recognized defect detection methods, including Ganomaly [10], KDAD [44], AnoViT [17], and DAAD [6], AESSIM [5], fAnoGan [9], DBISD [30], CAD [28], SCADN [7]. This comparative analysis allowed us to position our model on current standards in the field and to evaluate its performance in a quantifiable manner. To measure the effectiveness of these different methods, including ours, we opted for the use of the evaluation matrix of the area under the curve (AUC) of the receiver operating characteristics (ROC). This metric is widely recognized for its ability to provide a reliable and comprehensive evaluation of binary classification model performance, taking into account both the sensitivity, and specificity of the model. In addition to this, we also considered the F1-score and accuracy, which provide further insights into the model's precision and overall performance.

TABLE II. Comparison of pixel-level detection on the MVTec AD dataset.

Class	AnoGAN[8]	SMAI[26]	KDAD[44]	HaloAE[20]	GAP[33]	FCDD[40]	VTBA[36]	OUR
Bottle	86	86	96.3	91.9	93	97	95.1	94.1
Cable	78	92	82.4	87.6	94	90	92.6	87.1
Capsule	84	93	95.9	97.8	90	93	93.1	97.8
Hazelnut	87	97	94.6	97.8	84	95	98.2	97.7
Metal nut	76	92	86.4	85.2	91	94	91	90.6
Pill	87	92	89.6	91.5	93	81	92.6	98.6
Screw	80	96	96.0	99.0	96	86	97.7	97.6
Toothbrush	90	96	96.1	92.9	96	94	89.4	99.1
Transistor	80	85	76.5	87.5	100	88	85	87.9
Zipper	78	90	93.9	96.0	99	92	93.2	96.4
<i>Mean_{obj}</i>	82.6	91.9	90.7	92.7	93.6	91	92.8	94.7
Carpet	54	88	95.6	89.4	96	96	88.4	85.6
Grid	58	97	91.8	83.1	78	91	97.2	97.6
Leather	64	86	98.1	98.5	90	98	96.6	99.4
Tile	50	62	82.8	78.5	80	91	92.8	95
Wood	62	80	84.8	91.1	81	88	91.4	84
<i>Mean_{tex}</i>	57.6	82.6	90.6	88.1	85	92.8	93.3	92.3
Mean	74	89	90.7	91.2	91	92	93	93.9

TABLE III. Comparison of image-level results on the MVTec AD dataset.

Class	Ganomaly[10]	AnoVit[17]	KDAD[44]	DAAD[6]	DBISD[30]	OUR
Bottle	89.2	83	99.4	97.6	94	99.4
Cable	75.5	74	89.2	84.4	88	79
Capsule	73.2	73	80.5	76.7	85	82.7
Hazelnut	78.5	88	98.4	92.1	95	96.6
Metal nut	70.0	86	73.6	75.8	69	82.6
Pill	74.3	72	82.7	90.0	89	90.8
Screw	74.6	100	83.3	98.7	100	89.8
Toothbrush	65.3	74	92.2	99.2	100	99.7
Transistor	79.2	83	85.6	87.6	88	95.4
Zipper	74.5	73	93.2	85.9	91	98
<i>Mean_{obj}</i>	75.4	80.6	87.8	88.8	89.9	91.4
Carpet	69.9	50	79.3	86.6	91	57.2
Grid	70.8	52	78.0	95.7	94	94.7
Leather	84.2	85	95.1	86.2	95	100
Tile	79.4	89	91.6	88.2	80	99.4
Wood	83.4	95	94.3	98.2	94	96.2
<i>Mean_{tex}</i>	77.5	74.2	87.7	91	90.8	89.5
Mean	76.2	78	87.7	89.5	90.2	90.8

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (8)$$

where,

- TP: True Positives.
- TN: True Negatives.
- FP: False Positives.
- FN: False Negatives.

$$F1\text{-score} = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (9)$$

where,

$$Precision = \frac{TP}{TP + FP} \quad (10)$$

$$Recall = \frac{TP}{TP + FN} \quad (11)$$

Table II illustrates the performance achieved using the area under the ROC curve (AUC) for receiver operating characteristics at the pixel level. By comparing our model with some leading models in MVTEC AD, including two using transformers (HaloAE, VTBA) and others using convolutions, our model stands out markedly in five categories, with a lead ranging from 0.4% to 5.6%. It is worth noting that HaloAE utilizes data augmentation techniques, whereas our model does not employ any data augmentation. This remarkable superiority is also observed in the overall average of all object categories, where our model exceeds other methods by 1.1%. It is important to mention that both KDAD and FCDD are pre-trained models, while our model only uses pre-training for the first three blocks. Additionally, these categories represent 67% of the total data. Taking into account the overall average for all data categories, the performance of our method exceeds those of other compared methods by 0.9%.

Table III provides a comparison of the MVTEC AD dataset’s image-level detection findings, revealing that our model surpasses other models in five categories, with a lead ranging from 0.8% to 7.8%. In addition, the average of all categories for our model is higher than the total average of other approaches. Despite a negative impact observed in the ‘carpet’ category, these results highlight our model’s ability to effectively detect defects.

Table IV which presents the results obtained with the AUC of the ROC curve, comparing our model with state-of-the-art methods such as KDAD, AESSIM, CAD, fAnoGan, and SCADN on the constructed data class, as the class we created was not included in the original articles. The results reveal that our model surpasses KDAD by 6.8% in the button category. Combining the results of Tables II and IV, our model outperforms KDAD in 13 of the 16 categories at the pixel level and offers nearly the best results at the image level. It is important to note that KDAD uses the Teacher-Student mechanism, where the Teacher is pre-trained on ImageNet, while our model relies only on the initial pre-trained layers. Additionally, our model outperforms fAnoGan, CAD, AESSIM, and SCADN by 1% at the image level and exceeds AESSIM by 12.4% at the pixel level. These results highlight the effectiveness of our approach for defect detection and underscore the advantage of combining both local and global features, which enhances the accuracy of defect identification.

TABLE IV. Comparison results for the constructed class ‘Button’.

Model	pixel-level	image-level
OUR	97.4	99.7
AESSIM[5]	85	94
CAD[28]	-	95.8
fAnoGan[9]	-	98.5
SCADN[7]	-	91.3
KDAD[44]	90.6	99.5

Figures 5, 8 and 9 represent the evaluation of the model’s performance in terms of the visual localization of defects. In Figure 5, we focus on the constructed data class, while Figures 8 and 9 are concerned with the MVTEC AD dataset. Each line in the figures presents six columns: the first and fourth columns show the input image, the second and fifth columns display the segmentation mask, and the third and sixth columns represent the anomaly scores, where the red color indicates a high anomaly score. These representations demonstrate the proposed model’s ability to localize defects, whether their size is small, medium, or large. It is notable that the classes where the model excels in terms of score also demonstrate excellent visual localization of defects, as in the classes of toothbrush, leather, capsule and screw, even for very small sizes compared to other classes such as carpet, and wood, the model manages to provide accurate localization of anomalies, thus demonstrating the robustness of its approach.

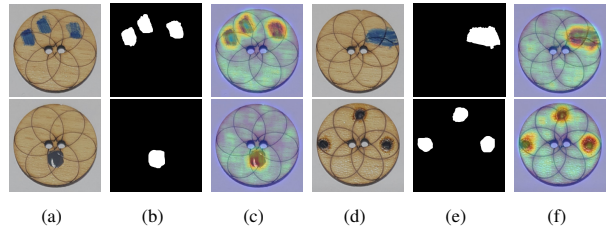


Figure 5. Qualitative results of our model on the ‘Button’ class. Rows a and d: Input images, b and e: Ground truth, c and f: Anomaly maps.

D. Ablation study

To demonstrate the effectiveness of our model and the impact of adding noise as well as skip connections, we conducted an ablation study. In the first case, we removed the noise, and in the second case, we eliminated the skip connections. Then, we compared these results with the full configuration. As shown in the Tables V, VI, VIII, VII and Figure 6 we evaluated the performance obtained under different configurations: without noise, without skip connections, and with both. The improvements brought by these elements are clearly visible, particularly in terms of accuracy, F1-score, and ROC AUC.

TABLE V. pixel-level Ablation results on certain classes of the MVTEC AD dataset: N noise, S skip.

class	Without N auc	Without S auc	With (S & N) auc
bottle	90	91.9	94.1
cable	75.8	78	79
carpet	59.5	85.3	87.1
grid	75.9	95	97.6
leather	97.5	97.9	99.6
toothbrush	98.4	99	99.7
tile	88.6	62.1	95
wood	78.3	74.7	84
zipper	91.1	94.8	96.4

TABLE VI. image-level Ablation results on certain classes of the MVTecAD dataset.

class	Without noise			Without skip			With (skip& noise)		
	f1-score	accuracy	auc	f1-score	accuracy	auc	f1-score	accuracy	auc
bottle	93	89	97.2	96	94	98.2	98	96	99.4
cable	76	61	54.8	76	61	62.7	77	71	79
carpet	87	77	45.3	87	77	57	88	79	57.2
grid	85	74	66.2	92	87	93.1	92	88	94.7
leather	85	75	82.9	91	87	93.4	100	100	100
toothbrush	89	86	97.8	97	95	98.6	98	98	99.7
tile	87	80	88.8	84	72	69.4	97	96	99.4
Wood	89	81	84.1	90	85	88.2	95	92	96.2
zipper	88	80	87.3	94	90	91.9	96	94	98

TABLE VII. Ablation results on MVTecAD dataset.

Model	pixel-level	image-level
With (noise&skip)	93.9	90.8
Without skip	90.1	87.6
Without noise	83.2	77.8

TABLE VIII. Ablation results on Button class.

Model	pixel-level	image-level
With (noise&skip)	97.4	99.4
Without skip	93	95
Without noise	94	99

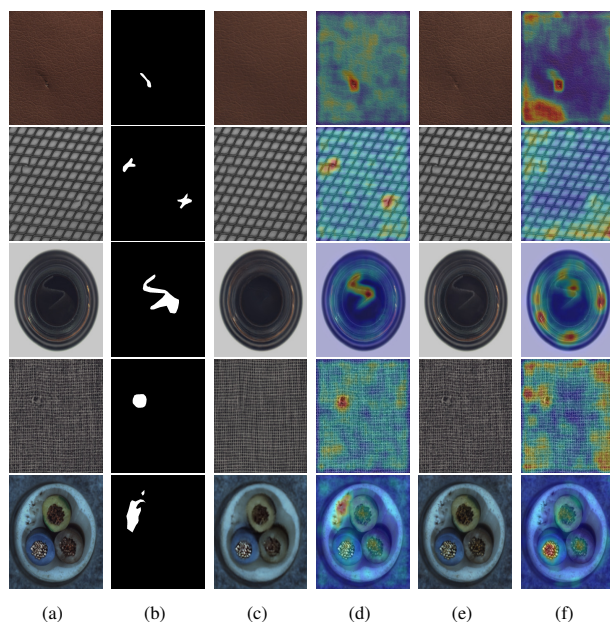


Figure 6. Qualitative Comparison with and without Noise on Certain Classes of the MVTec AD Dataset: a) Input Image, b) Ground Truth, c) Reconstructed Image with Noise, d) Heatmap with Noise, e) Reconstructed Image without Noise, f) Heatmap without Noise.

5. REAL-TIME DEFECT DETECTION INTEGRATION

In an automated industrial production line, the inference time shown in Image 7 indicates that our model can be effectively integrated into a real-time defect detection system using the following approach:

- High-resolution IoT cameras equipped with smart sensors are installed along the production line to capture real-time images of products at various stages. These cameras are connected to an industrial network via low-latency transmission protocols such as LAN, Profinet, or EtherCAT. These networks are designed to ensure fast and reliable transmission of data to the analysis system with minimal delay, which is crucial for maintaining production efficiency.
- Once the images are transmitted to the analysis system, they undergo necessary pre-processing steps, such as resizing and normalization. This ensures that the images conform to the input requirements of the defect detection model, allowing for efficient and accurate analysis.

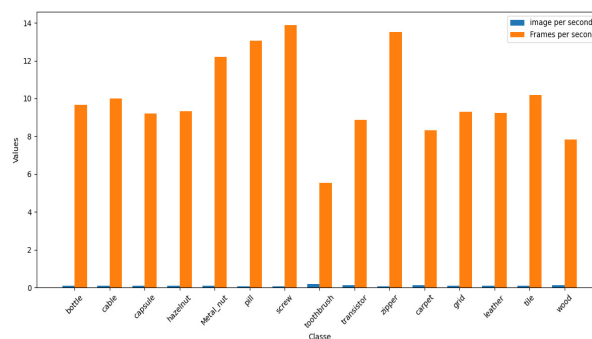


Figure 7. Inference time of ours model in the MVTec AD dataset.

- The pre-processed images are then fed into our model, where defects or anomalies are detected in real time. Thanks to the model's computational efficiency, as outlined in our earlier response regarding inference time, the system can operate at near real-time speeds, ensuring that defects are identified as products

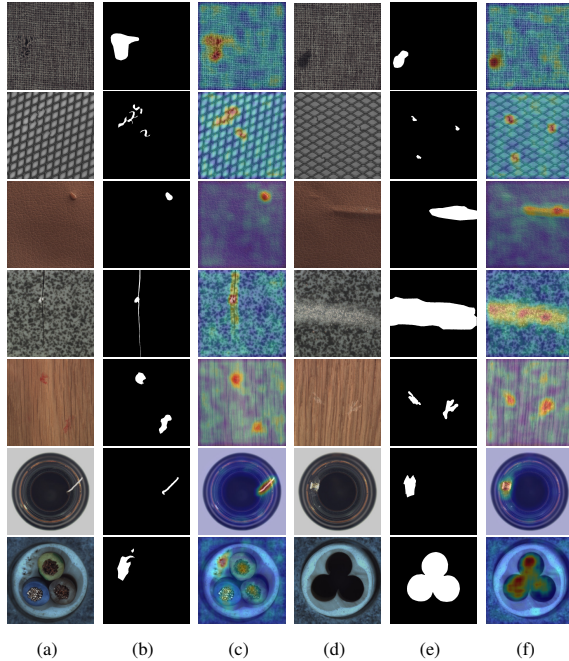


Figure 8. Qualitative results of our model on 7 out of the 15 classes in the MVTEC AD database. Rows a and d: Input images, b and e: Ground truth, c and f: Anomaly maps.

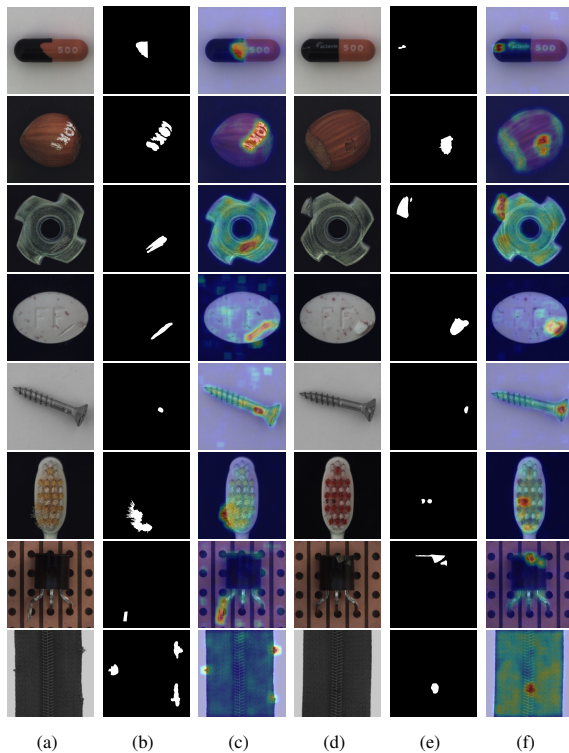


Figure 9. Qualitative results of our model on 8 out of the 15 classes in the MVTEC AD database. Rows a and d: Input images, b, and e: Ground truth, c, and f: Anomaly maps.

pass through the production line.

- The results of the analysis are displayed in real-time on an intuitive user interface, enabling operators to continuously monitor product quality. This interface can be customized to show detailed information about detected defects, including their location and severity.

6. LIMITATIONS OF OURS MODEL

The main challenge with this model lies in the optimal management of noise. If the noise level is too low or completely removed, the model may reconstruct the defects present in the image, making their localization impossible. Conversely, an excessive increase in noise to prevent the reconstruction of defects leads to poor overall image quality. This degradation directly affects the accuracy of defect localization and detection, thereby reducing the overall effectiveness of the model, as seen in cases of anomalies on surfaces such as carpets, cables, or metal nuts.

7. CONCLUSION

In this work, we have developed an innovative architecture that combines convolutional neural networks (CNNs) and transformers to leverage their respective strengths in extracting local and global features. Our encoder-decoder architecture is distinguished by the integration of CNN blocks, which are pre-trained to capture fine and local details in the early layers of the encoder. This approach is complemented by the use of transformers in the final layers to capture and integrate information on a broader scale. The skip connections from the encoder to the decoder especially the first two, which are reinforced by multi-head attention and spatial reduction, followed by a convolutional operation play a crucial role in effectively weighting the features relevant to the decoding task. Moreover, the introduction of random Gaussian noise upstream of the image contributes to preventing the reconstruction of defects, representing a significant step toward model robustness. Establishing a specific data class marks a notable advancement in our research methodology.

Our future work will mainly focus on improving the model's performance in classes that currently have a negative impact on our results. This approach will involve a thorough analysis of these specific categories to identify challenges and obstacles that hinder effective processing. We will consider integrating new deep learning techniques and optimizing the architecture to refine the model's ability to handle more complex or unconventional cases. Additionally, we will explore the effectiveness of different types of noise and regularization techniques to further enhance the model's ability to generalize and avoid overfitting, particularly in scenarios where data is limited or highly specific. The ultimate goal of this future work will be to improve the robustness and accuracy of the model, making it more effective and adaptable to various practical applications.

REFERENCES

- [1] M. Javaid, A. Haleem, R. P. Singh, S. Rab, and R. Suman, "Exploring impact and features of machine vision for progressive industry 4.0 culture," *Sensors International*, vol. 3, p. 100132, 2022.
- [2] F. K. Konstantinidis, S. G. Mouroutsos, and A. Gasteratos, "The role of machine vision in industry 4.0: an automotive manufacturing perspective," in *2021 IEEE International Conference on Imaging Systems and Techniques (IST)*, pp. 1–6, 2021.
- [3] M. Javaid, A. Haleem, R. P. Singh, and R. Suman, "Artificial intelligence applications for industry 4.0: A literature-based study," *Journal of Industrial Integration and Management*, vol. 7, pp. 83–111, 2022.
- [4] R. Rai, M. K. Tiwari, D. Ivanov, and A. Dolgui, "Machine learning in manufacturing and industry 4.0 applications," *International Journal of Production Research*, vol. 59, no. 16, pp. 4773–4778, 2021.
- [5] P. Bergmann, S. Löwe, M. Fauser, D. Sattlegger, and C. Steger, "Improving unsupervised defect segmentation by applying structural similarity to autoencoders," *arXiv preprint arXiv:1807.02011*, 2018.
- [6] J. Hou, Y. Zhang, Q. Zhong, D. Xie, S. Pu, and H. Zhou, "Divide-and-assemble: Learning block-wise memory for unsupervised anomaly detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*. IEEE, 2021, pp. 8791–8800.
- [7] X. Yan, H. Zhang, X. Xu, X. Hu, and P.-A. Heng, "Learning semantic context from normal samples for unsupervised anomaly detection," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 35, 2021, pp. 3110–3118.
- [8] T. Schlegl, P. Seeböck, S. M. Waldstein, U. Schmidt-Erfurth, and G. Langs, "Unsupervised anomaly detection with generative adversarial networks to guide marker discovery," in *International conference on information processing in medical imaging*, pp. 146–157, 2017.
- [9] T. Schlegl, P. Seeböck, S. M. Waldstein, G. Langs, and U. Schmidt-Erfurth, "f-anogan: Fast unsupervised anomaly detection with generative adversarial networks," *Medical image analysis*, vol. 54, pp. 30–44, 2019.
- [10] S. Akcay, A. Atapour-Abarghouei, and T. P. Breckon, "Ganomaly: Semi-supervised anomaly detection via adversarial training," in *Computer Vision—ACCV 2018: 14th Asian Conference on Computer Vision, Perth, Australia, December 2–6, 2018, Revised Selected Papers, Part III*, vol. 14, pp. 622–637, 2019.
- [11] D. Zimmerer, S. A. Kohl, J. Petersen, F. Isensee, and K. H. Maier-Hein, "Context-encoding variational autoencoder for unsupervised anomaly detection," *arXiv preprint arXiv:1812.05941*, 2018.
- [12] Y. Lu and P. Xu, "Anomaly detection for skin disease images using variational autoencoder," *arXiv preprint arXiv:1807.01349*, 2018.
- [13] A. A. Pol, V. Berger, C. Germain, G. Cerminara, and M. Pierini, "Anomaly detection with conditional variational autoencoders," in *2019 18th IEEE international conference on machine learning and applications (ICMLA)*, pp. 1651–1657, 2019.
- [14] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [15] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [16] P. Mishra, R. Verk, D. Fornasier, C. Piciarelli, and G. L. Foresti, "Vt-adl: A vision transformer network for image anomaly detection and localization," in *30th International Symposium on Industrial Electronics (ISIE)*, pp. 01–06, 2021.
- [17] Y. Lee and P. Kang, "Anovit: Unsupervised anomaly detection and localization with vision transformer-based encoder-decoder," *IEEE Access*, vol. 10, p. 46717–46724, 2022.
- [18] Y. Yang, "Self-attention autoencoder for anomaly segmentation," 2021.
- [19] H. Wu, B. Xiao, N. Codella, M. Liu, X. Dai, L. Yuan, and L. Zhang, "Cvt: Introducing convolutions to vision transformers," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 22–31.
- [20] E. Mathian, H. Liu, L. Fernandez-Cuesta, D. Samaras, M. Foll, and L. Chen, "Haloae: An halonet based local transformer auto encoder for anomaly detection and localization," *arXiv preprint arXiv:2208.03486*, 2022.
- [21] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5–9, 2015, proceedings, part III*, vol. 18, pp. 234–241, 2015.
- [22] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, "Pyramid vision transformer: A versatile backbone for dense prediction without convolutions," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 568–578.
- [23] J. Yang, Y. Shi, and Z. Qi, "Dfr: Deep feature reconstruction for unsupervised anomaly segmentation," *arXiv preprint arXiv:2012.07122*, 2020.
- [24] V. Zavrtnik, M. Kristan, and D. Skočaj, "Reconstruction by inpainting for visual anomaly detection," *Pattern Recognition*, vol. 112, p. 107706, 2021.
- [25] K. Zhou, Y. Xiao, J. Yang, J. Cheng, W. Liu, W. Luo, Z. Gu, J. Liu, and S. Gao, "Encoding structure-texture relation with p-net for anomaly detection in retinal images," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XX*, vol. 16, pp. 360–377, 2020.
- [26] Z. Li, N. Li, K. Jiang, Z. Ma, X. Wei, X. Hong, , and Y. Gong, "Superpixel masking and inpainting for self-supervised anomaly detection," in *Bmvc*, 2020.
- [27] R. Jiang, Y. Xue, and D. Zou, "Interpretability-aware industrial anomaly detection using autoencoders," *IEEE Access*, vol. 11, pp. 60 490–60 500, 2023.
- [28] W. Li, J. Zhan, J. Wang, B. Xia, B.-B. Gao, J. Liu, C. Wang, and F. Zheng, "Towards continual adaptation in industrial anomaly detection," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 2871–2880.



- [29] L. Zhao, Y. Chai, Q. Zhang, and H. R. Karimi, "Self-supervised anomaly detection based on foreground enhancement and auto-encoder reconstruction," *Signal, Image and Video Processing*, vol. 18, no. 1, pp. 343–350, 2024.
- [30] X. Li, Y. Zheng, B. Chen, and E. Zheng, "Dual attention-based industrial surface defect detection with consistency loss," *Sensors*, vol. 22, no. 14, p. 5141, 2022.
- [31] H. Zhang, S. Wang, S. Lu, L. Yao, and Y. Hu, "Attention-gate-based u-shaped reconstruction network (agur-net) for color-patterned fabric defect detection," *Textile Research Journal*, vol. 93, no. 15–16, pp. 3459–3477, 2023.
- [32] R. Chen, G. Xie, J. Liu, J. Wang, Z. Luo, J. Wang, , and F. Zheng, "Easynet: An easy network for 3d industrial anomaly detection," in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 7038–7046.
- [33] S. Wang, L. Wu, L. Cui, and Y. Shen, "Glancing at the patch: Anomaly localization with global and local feature comparison," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 254–263.
- [34] A. De Nardin, P. Mishra, G. L. Foresti, and C. Piciarelli, "Masked transformer for image anomaly localization," *International Journal of Neural Systems*, vol. 32, no. 7, p. 2250030, 2022.
- [35] K. Y. You Zhiyuan, W. Luo, L. Cui, Y. Zheng, and X. Le, "Adtr: anomaly detection transformer with feature reconstruction," in *International Conference on Neural Information Processing*, pp. 298–310, 2022.
- [36] Q. Yang and R. Guo, "An unsupervised method for industrial image anomaly detection with vision transformer-based autoencoder," *Sensors*, vol. 24, no. 8, p. 2440, 2024.
- [37] H. Shang, C. Sun, J. Liu, X. Chen, and R. Yan, "Defect-aware transformer network for intelligent visual surface defect detection," *Advanced Engineering Informatics*, vol. 55, p. 101882, 2023.
- [38] J. Yi and S. Yoon, "Patch svdd: Patch-level svdd for anomaly detection and segmentation," in *Proceedings of the Asian conference on computer vision*, 2020.
- [39] F. V. Massoli, F. Falchi, A. Kantarci, Şeymanur Akti, H. K. Ekenel, and G. Amato, "Mocca: Multilayer one-class classification for anomaly detection," *IEEE Transactions on neural networks and learning systems*, vol. 33, no. 6, pp. 2313–2323, 2021.
- [40] P. Liznerski, L. Ruff, R. A. Vandermeulen, B. J. Franks, M. Kloft, and K.-R. Müller, "Explainable deep one-class classification," *arXiv preprint arXiv:2007.01760*, 2020.
- [41] J. Yoo, L. Zhao, and L. Akoglu, "End-to-end augmentation hyperparameter tuning for self-supervised anomaly detection," *arXiv preprint arXiv:2306.12033*, 2023.
- [42] O. Rippel, A. Chavan, C. Lei, and D. Merhof, "Transfer learning gaussian anomaly detection by fine-tuning representations," *arXiv preprint arXiv:2108.04116*, 2021.
- [43] O. Rippel, P. Mertens, and D. Merhof, "Modeling the distribution of normal data in pre-trained deep features for anomaly detection," in *2020 25th International Conference on Pattern Recognition (ICPR)*, pp. 6726–6733, 2021.
- [44] M. Salehi, N. Sadjadi, S. Baselizadeh, M. H. Rohban, and H. R. Rabiee, "Multiresolution knowledge distillation for anomaly detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 14902–14912.
- [45] Q. Wu, H. Li, C. Tian, L. Wen, and X. Li, "Aekd: Unsupervised auto-encoder knowledge distillation for industrial anomaly detection," *Journal of Manufacturing Systems*, vol. 73, pp. 159–169, 2024.
- [46] W. Xue, L. Zhang, X. Mou, and A. C. Bovik, "Gradient magnitude similarity deviation: A highly efficient perceptual image quality index," *IEEE transactions on image processing*, vol. 23, no. 2, pp. 684–695, 2013.
- [47] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger, "Mvtec ad—a comprehensive real-world dataset for unsupervised anomaly detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9592–9600.