

Bridge Feature Matching and Cross-Modal Alignment with Mutual-filtering for Zero-shot Anomaly Detection

Yuhu Bai^{1*} Jiangning Zhang^{1,2*} Yunkang Cao³ Guangyuan Lu¹
Qingdong He² Xiangtai Li⁴ Guanzhong Tian^{1†}

¹Zhejiang University ²YouTu Lab, Tencent ³Hunan University ⁴Peking University

Code: <https://github.com/ybai111/FiSeCLIP>

Abstract

With the advent of vision-language models (e.g., CLIP) in zero- and few-shot settings, CLIP has been widely applied to zero-shot anomaly detection (ZSAD) in recent research, where the rare classes are essential and expected in many applications. This study introduces *FiSeCLIP* for ZSAD with training-free CLIP, combining the feature matching with the cross-modal alignment. Testing with the entire dataset is impractical, while batch-based testing better aligns with real industrial needs, and images within a batch can serve as mutual reference points. Accordingly, *FiSeCLIP* utilizes other images in the same batch as reference information for the current image. However, the lack of labels for these references can introduce ambiguity, we apply text information to filter out noisy features. In addition, we further explore CLIP’s inherent potential to restore its local semantic correlation, adapting it for fine-grained anomaly detection tasks to enable a more accurate filtering process. Our approach exhibits superior performance for both anomaly classification and segmentation on anomaly detection benchmarks, building a stronger baseline for the direction, e.g., on MVTec-AD, *FiSeCLIP* outperforms the SOTA AdaCLIP by +4.6% \uparrow +5.7% \uparrow in segmentation metrics AU-ROC/ F_1 -max.

1. Introduction

Anomaly detection (AD) has found extensive applications in various domains [6, 29, 30, 44], where it plays a crucial role. AD aims to identify if a sample has any anomalies and to pinpoint the anomalous locations. Previous anomaly detection [12, 22, 35, 38, 41] approaches mainly follow an

*Equal contribution.

†Corresponding author. This work is supported in part by the National Natural Science Foundation of China under Grant 62303405, in part by Ningbo Natural Science Foundation Project under Grant 2023J400, and in part by Open Research Fund Program of Beijing National Research Center for Information Science and Technology.

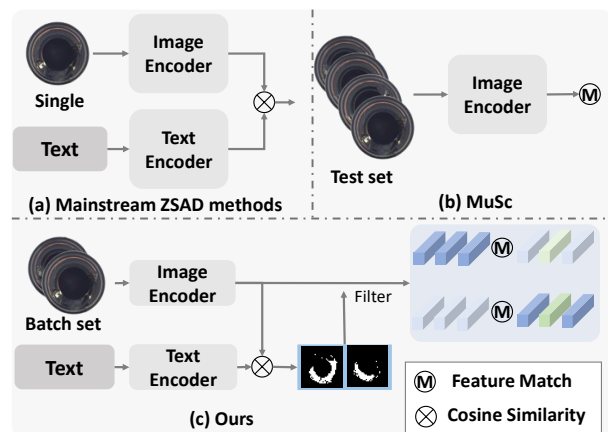


Figure 1. Compared with (a) mainstream ZSAD methods [9, 19] and (b) MuSc [26] method. Our method uses a batch of images, whereas previous methods rely on either a single image or the entire dataset.

unsupervised paradigm, yet they still rely on a substantial amount of normal samples for training. However, requiring substantial computational resources, these methods also lack strong generalization abilities. Zero-shot anomaly detection (ZSAD) has emerged [19, 24], allowing direct inference without any target domain data, and has recently garnered significant attention from researchers.

The advancement of vision-language pre-training [34] and their strong zero-shot generalization abilities have led many recent studies to incorporate them into zero-shot anomaly detection. In particular, WinCLIP [19] and APRIL-GAN [9] initially employed CLIP [34] in anomaly detection, leveraging the similarity between image and text features to estimate anomaly probability. As CLIP [34] was developed for classification tasks, some works [10, 14, 47] observe that the original $q - k$ attention disrupts local semantics and limits fine-grained anomaly segmentation, where q, k, v represent the *query, key, value* embeddings. They introduced $v - v$ attention to strengthen local se-

semantic representation. Although improvements in the self-attention mechanism have led to performance gains, local semantic representation remains insufficient, and modifications across multiple layers have somewhat compromised the model’s generalization ability. In addition to improvements on the visual level, studies like AnomalyCLIP [47], AdaCLIP [7] and Filo [14] have also designed learnable text prompts, aiming to enable the model to understand the concepts of ‘normal’ and ‘abnormal’. Additionally, MuSc [26] finds that over 95% of the pixels in the test set of MVTec AD dataset [4] and VisA dataset [48] are normal, so it utilizes the entire test dataset for feature matching in anomaly detection. However, There are certain problems associated with the aforementioned methods. First, in practical production environments, testing is usually conducted on batches of samples, while applying the entire test dataset in real-world industrial scenarios is impractical and requires extensive computational resources. Moreover, the posterior optimization in MuSc [26] requires prior knowledge of the test dataset distribution before testing. Second, most of ZSAD methods involve fine-tuning, a setup that may risk data leakage.

To tackle the abovementioned issues, we propose training-free FiSeCLIP, which combines mainstream zero-shot anomaly detection principles with the idea of feature matching. *During testing, we process images in batches, leveraging cross-reference cues within the batch to extract limited reference information, thus requiring simultaneous input of all batch images to the model.* As shown in Fig 1, we will apply a strategy similar to MuSc [26] and PatchCore [35], denoted FiCLIP-AD. As input images are unlabeled test images, this can lead to noise interference. In contrast to MuSc[26], we propose a filtering strategy to filter out anomalous features from the unknown label reference information. The *filtering mask* is dynamically constructed via cross-modal similarity alignment between CLIP’s visual embeddings and their corresponding textual embeddings, named SeCLIP-AD, to optimize feature-matching outcomes. Specifically, we find that intermediate-layer attention provides stronger local semantic representation, while the residual connection and the feed-forward network can degrade fine-grained segmentation performance [20]. To retain CLIP’s generalization ability, we replace only the final layer’s attention with intermediate-layer attention and remove the feed-forward network. Furthermore, by integrating the results of FiCLIP-AD and SeCLIP-AD, the fused outcomes enable iterative refinement of the *filtering mask*. This mutual optimization loop enhances the accuracy of FiCLIP-AD’s anomaly score map through continuous feedback, where the refined *filtering mask* guides more precise feature matching and noise suppression. In summary, our contributions can be outlined as follows:

- We propose a novel method, FiSeCLIP, which integrates feature matching with mainstream ZSAD techniques. By leveraging unlabeled test images within the same batch as mutual references, FiSeCLIP introduces a mutual-filtering strategy to dynamically filter out anomalous features and mitigate noise interference during feature matching.
- We begin by using textual information to perform initial filtering of noisy features, followed by further refinement with feature matching results. We propose enhancing CLIP’s semantic coherence to improve the alignment of fine-grained visual and textual features.
- We perform comprehensive experiments, achieving state-of-the-art performance on widely used benchmarks, e.g., we achieved improvements of +4.6%↑/+5.7%↑ and +6.1%↑/+4.8%↑ in AU-ROC/ F_1 -max for anomaly segmentation and classification, respectively.

2. Related Work

2.1. Industrial anomaly detection

Most industrial anomaly detection works can be classified into three categories: 1) reconstruction-based methods; 2) synthesizing-based methods and 3) embedding-based methods. **Reconstruction-based methods** [16, 18, 28, 38, 40, 43] posit that anomalous regions are difficult to reconstruct correctly. Thus, reconstruction error serves as an indicator for locating anomalous areas. **Synthesizing-based methods** [8, 25, 42, 45] generate artificial anomalies to provide supervision signals for training. **Embedding-based methods** [2, 3, 11, 15, 31, 35] extract features using a pre-trained network, then establish a memory bank to detect anomalous patches through feature matching. However, these methods require a large number of samples and have limited generalization capability.

2.2. Vision-language models for ZSAD

Vision-language pre-training models. Vision-language pre-training has emerged as a powerful approach for visual representation learning, with CLIP [34] standing out due to its impressive generalization capabilities. Pre-trained on a large dataset of images sourced from the Web, CLIP aligns images and natural language through two separate encoders, typically based on architectures like ResNet [17], ViT [13], or their enhanced versions. This design allows CLIP to be easily adapted to various downstream classification tasks using prompts. Although originally intended for classification, CLIP has been adapted for fine-grained segmentation [20, 21, 27, 37, 39] without any training. These works improve local feature awareness with a modified attention layer, addressing the dominance of global patches, and boosting segmentation accuracy and semantic coherence. The anomaly detection task is closely related to

the the fine-grained segmentation task, and similar research ideas have been explored and extended to anomaly detection in later studies.

Zero-shot anomaly detection. WinCLIP [19] pioneers the application of CLIP in zero-shot anomaly detection. Subsequently, APRIL-GAN [9] improves upon it by incorporating a linear layer for fine-tuning. Some works [10, 14, 47] modify the $q - k$ attention to $v - v$ attention to enhance CLIP’s local semantics. Moreover, several methods [7, 14, 33, 47] employ learnable text prompts to boost zero-shot generalization, with the goal of enabling the model to autonomously identify the concepts of ‘anomaly’ and ‘normal’. However, the limited alignment of CLIP’s visual and textual features in fine-grained tasks severely affects ZSAD performance. Furthermore, the need for fine-tuning in many prior approaches raises the risk of data leakage. Recently, MuSc [26] has demonstrated outstanding performance by taking advantage of the entire test set with the idea of matching features. However, using the entire data set is impractical for real-world anomaly detection and requires substantial memory resources.

3. Method

This section presents our proposed FiSeCLIP, divided into two components: FiCLIP-AD and SeCLIP-AD. FiCLIP-AD uses CLIP to extract features from unlabeled images within a batch, with each image as a reference for the others. However, unlabeled references may introduce noise interference. To address this, we propose SeCLIP-AD to first generate a prior filtering mask that suppresses anomalous features, ensuring cleaner reference information for subsequent feature matching. SeCLIP-AD investigates a training-free strategy that harnesses CLIP’s intrinsic capabilities to enhance local semantic self-correlation among patches. The anomaly score A_{Se} is obtained by calculating the similarity between visual and text features. We seamlessly integrate FiCLIP-AD and SeCLIP-AD predictions and iteratively feed them back to FiCLIP-AD, significantly enhancing anomaly score map accuracy while suppressing false positives.

3.1. Task setting.

In this study, we perform simultaneous testing on images within a batch $D = \{I_u\}, u = 1, \dots, B$, using them as mutual references, a method we refer to as batch zero-shot anomaly detection. In contrast, mainstream zero-shot anomaly detection [7, 9, 19, 47] requires only a single image $D = \{I_u\}, u = 1$ for inference. The recent MuSc [26] requires the entire test dataset $D = \{I_u\}, u = 1, \dots, n$ to be processed simultaneously for inference. MuSc-2 [26] and Dual Image Enhanced CLIP [46] adopt two images from the test set for evaluation. The remaining methods [7, 9, 19, 47] all take a single image as input and output the results. For

simplicity, we will use two images (I_u and I_v) as examples to introduce our method in the following sections.

3.2. SeCLIP-AD

3.2.1. Semantic correlation recovering

CLIP is trained for classification tasks with more attention to global semantic information. Meanwhile, anomaly classification and segmentation are required to understand local semantics. Thus, we need to recover its semantic correlation for the local patterns.

A ViT-based CLIP vision encoder comprises a series of attention blocks. Each block yield the feature representation $X_i \in \mathbb{R}^{B \times (HW+1) \times D}$, where i represents the layer index, B is the batch size, D is the dimension, HW is the number of local patch tokens, and the other one is $[CLS]$ token. For simplicity, a self-attention-based ViT encoder is described as follows:

$$[q_i, k_i, v_i] = \text{Proj}_{qkv}(\text{LN}(X_{i-1})), \quad (1)$$

$$\text{Attn}_i = \text{Softmax}\left(\frac{q_i^T k_i}{\sqrt{d_k}}\right), \quad (2)$$

$$X_i = X_{i-1} + X_{\text{attn}} = X_{i-1} + \text{Proj}(\text{Attn}_i \cdot v_i), \quad (3)$$

$$X_i = X_i + \text{FFN}(\text{LN}(X_i)), \quad (4)$$

where Proj denotes a projection layer, LN denotes layer normalization, and FFN represents a feed-forward network. q , k , and v represent the query, key, and value embeddings, respectively. i is the layer index, and $i - 1$ denotes the output from the previous layer.

To retain CLIP’s generalization capacity, we restrict changes to the last layer. Visualizing the CLIP’s attention map, as shown in Fig.3, it can be found that the middle layer’s attention map has better local semantics. As a result, the intermediate layer’s attention $\text{Attn}_{\text{inter}}$ is substituted for that of the final layer attention map Attn_{-1} :

$$\text{Attn}_{-1} = \text{Attn}_{\text{inter}}, X_{\text{attn}} = \text{Proj}(\text{Attn}_{\text{inter}} \cdot v_{-1}). \quad (5)$$

What’s more, ClearCLIP [20] finds residual connection significantly degrade performance on dense segmentation tasks, with the feed-forward module having a negligible impact. In this paper, we adopt a similar approach. Specifically, the output $X_{\text{final}} \in \mathbb{R}^{B \times (HW+1) \times D}$ of the final block layer can be defined as $X_{\text{final}} = X_{\text{attn}}$.

3.2.2. CLIP for Zero-shot Anomaly Detection

To leverage CLIP’s multi-modal capability for zero-shot anomaly detection, a naive idea is to compute the similarity between text features and visual features.

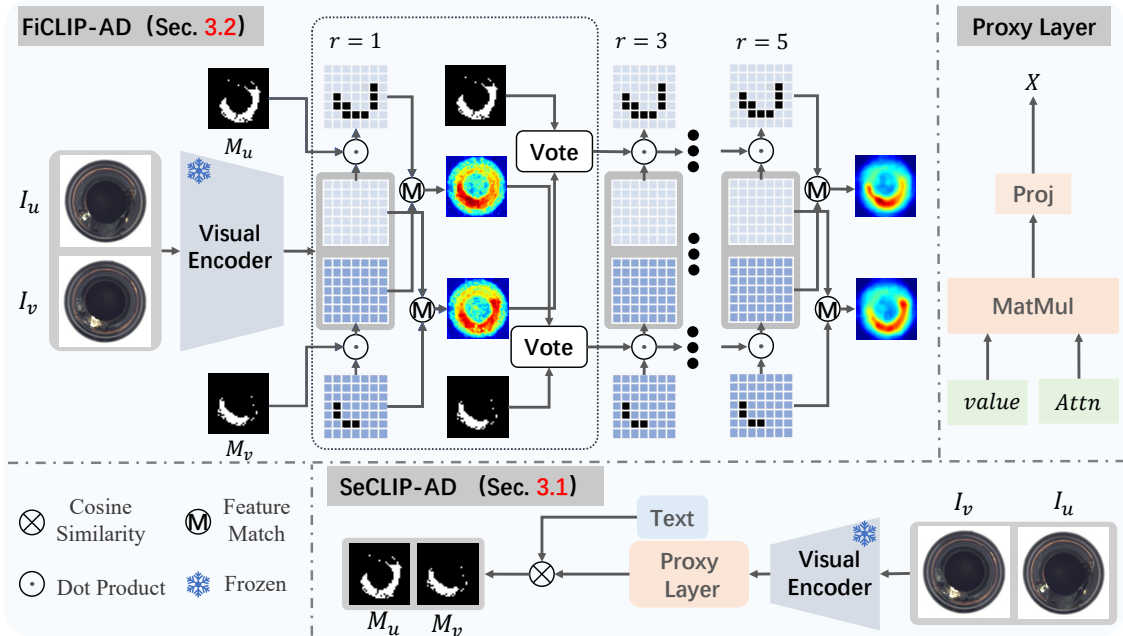


Figure 2. **Overview of our FiSeCLIP.** It consists of two parts: SeCLIP-AD (Sec. 3.2) and FiCLIP-AD (Sec. 3.3). 1) For two images in the same batch, SeCLIP-AD calculates the similarity between text and visual features to generate the corresponding anomaly score and mask, which are then used by the FiCLIP-AD module. The Proxy Layer is used to restore semantic correlation among patches, with attention weights derived from the intermediate layer. 2) For FiCLIP-AD, two images serve as inputs for feature extraction and aggregation with varying neighbor counts ($r = 1, 3, 5$). The anomaly mask from the previous step and the mask generated by SeCLIP-AD are refined through mutual voting, and the final mask is used to filter noisy features for feature matching. Different shades of blue represent patch tokens from different images.

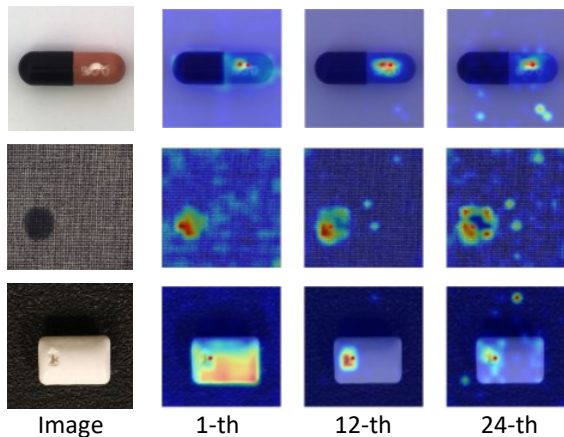


Figure 3. **Attention maps from different layers of CLIP**, with red dots indicating the selected patch positions in the image, which shows 12-th attention map has the best local semantics.

Anomaly detection involves the concepts of normal and abnormal, so the text prompts can be simply designed as ‘a [d] photo of [s] [c]’, where [s] represents the concepts of normality and abnormality, such as ‘normal’ and ‘abnor-

mal,’ or ‘flawless’ and ‘damaged.’ [d] offers additional information about the image, such as ‘rotated’ or ‘cropped,’ and [c] indicates the image’s category. Two types of text prompts, normal and abnormal, are generated. We utilize CLIP’s text encoder to extract text features, and the final text features $F_t \in \mathbb{R}^{2 \times C}$ can be obtained by averaging the features of normal and abnormal prompts separately, where C denotes the dimension of the features.

For anomaly classification, the visual global features $F_c \in \mathbb{R}^{B \times C}$ of original CLIP is utilized, and the probability of an object $CLS_{prob} \in \mathbb{R}^{B \times 2}$ being classified as normal or abnormal can be expressed as:

$$CLS_{prob} = \text{Softmax}(F_c \cdot F_t). \quad (6)$$

The anomaly classification score $CLS_{score} \in \mathbb{R}^B$ refers to the probability of the category being abnormal.

For anomaly segmentation, we remove the [CLS] token and map the dimension to C , ultimately obtaining the local patch features $F_s \in \mathbb{R}^{B \times HW \times C}$. The anomaly segmentation probability $SEG_{prob} \in \mathbb{R}^{B \times H \times W \times 2}$ can be described as:

$$SEG_{prob} = \text{Softmax}(F_s \cdot F_t), \quad (7)$$

and the anomaly segmentation score $Seg_{score} \in \mathbb{R}^{B \times HW}$ refers to the probability of the patch being abnormal.

3.3. FiCLIP-AD

3.3.1. Feature extraction

We use the original CLIP to simultaneously extract the features of B unlabeled test images $D = \{I_u\}, u = 1, \dots, B$. For the ViT-based CLIP encoder with L layers, the multi-stage patch tokens $F_u^i \in \mathbb{R}^{HW \times D}$ are employed, where $i \in \{0, 1, \dots, L\}$ indicates the layer of CLIP visual encoder. By aggregating neighboring features, we enhance the detection accuracy for anomalies across various sizes. Specifically, the patch tokens $F_u^i \in \mathbb{R}^{HW \times D}$ can be reshaped into $H \times W \times D$. Average pooling is employed to aggregate the patch token features in an $r \times r$ neighborhood of the current position to get the aggregated token $F_u^{i,r} \in \mathbb{R}^{H \times W \times D}$ following [26, 35], and reshape it to $F_u^{i,r} \in \mathbb{R}^{HW \times D}$.

3.3.2. Noisy features filtering

For the B images in the same batch size, the corresponding features are extracted using the method outlined in Sec.3.3.1. Calculate the distance between each patch token of the visual features $F_u^{i,r} \in \mathbb{R}^{HW \times D}$ and $F_v^{i,r} \in \mathbb{R}^{HW \times D}$, $v \neq u$, taking the minimum distance as the anomaly score $a_{u,v}^{i,r} \in \mathbb{R}^{HW}$:

$$a_{u,v}^{i,r} = \min \|F_u^{i,r} - F_v^{i,r}\|_2, \quad (8)$$

In applying feature-match methods[35, 36] for anomaly detection, we aim for the reference images to contain no patches similar to the anomalous patches in the inference image. However, when applying the idea to zero-shot anomaly detection, the uncertainty of whether the reference features are normal can introduce some interference. As shown in Fig.4, the similarity between anomalous patches in the reference and inference images results in a decreased anomaly score for the inference image, which may lead to a higher rate of false negatives in anomaly detection. Therefore, we propose filtering out the anomalous features in the reference features $F_v^{i,r}$ and using the filtered features as references to calculate the anomaly score. The initial anomaly mask $M \in \mathbb{R}^{HW}$ can be obtained in Eq.7:

$$M = \begin{cases} True, & \text{if } P_a > \lambda \cdot P_n \\ False, & \text{otherwise} \end{cases} \quad (9)$$

where P_a and P_n represent the abnormal and normal probabilities in Eq.7. Consequently, the ‘normal’ features can be defined as:

$$\hat{F}_v^{i,r} = F_v^{i,r}[\bar{M}], \quad (10)$$

where $\hat{F}_v^{i,r} \in \mathbb{R}^{N' \times D}$, and N' denotes the length of normal patches. The anomaly score is further described as follows:

$$\hat{a}_{u,v}^{i,r} = \min \|F_u^{i,r} - \hat{F}_v^{i,r}\|_2, \quad (11)$$

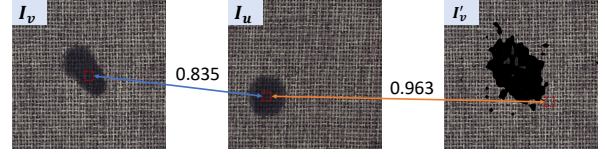


Figure 4. An example of mutual-filtering. I_u denotes the inference image, I_v represents the reference image without filtering, and I'_v represents the reference image with filtering.

Algorithm 1: Noise mutual filtering

Input: Starting masks M_u and M_v , inference features $F_u^{i,r}$, reference features $F_v^{i,r}$, stage index i , aggregation neighborhood r , initial anomaly score \hat{a} .

Output: Anomaly score \hat{a}_u, \hat{a}_v .

Algorithm:

$\hat{a}_u, \hat{a}_v \leftarrow empty$

for $r \in 1, 3, 5$ **do**

for $i \in 6, 12, 18, 24$ **do**

$\hat{F}_u^{i,r} \leftarrow F_u^{i,r}[\bar{M}_u], \hat{F}_v^{i,r} \leftarrow F_v^{i,r}[\bar{M}_v]$

$\hat{a}_u^{i,r} \leftarrow \min \|F_u^{i,r} - \hat{F}_v^{i,r}\|_2$

$\hat{a}_v^{i,r} \leftarrow \min \|F_v^{i,r} - \hat{F}_u^{i,r}\|_2$

end

$\hat{a}_u^r \leftarrow \frac{1}{4} \sum \hat{a}_u^{i,r}, \hat{a}_v^r \leftarrow \frac{1}{4} \sum \hat{a}_v^{i,r}$

$\hat{a}_u \leftarrow \text{mean}(\hat{a} + \hat{a}_u^r), \hat{a}_v \leftarrow \text{mean}(\hat{a}_v + \hat{a}_v^r)$

$M_u^{inter} \leftarrow \text{binary}(\hat{a}_u^r), M_v^{inter} \leftarrow \text{binary}(\hat{a}_v^r)$

$M_u \leftarrow \text{vote}(M_u, M_u^{inter})$

$M_v \leftarrow \text{vote}(M_v, M_v^{inter})$

end

3.3.3. Mutual-filtering mask

Moreover, Eq.11 allows one to obtain multiple anomaly scores. Therefore, we propose using the anomaly score to refine the anomaly mask M . Specifically, the four-stage patch tokens are utilized to generate the average anomaly score:

$$\hat{a}_v^r = \frac{1}{m} \sum_{i \in \{0,1,\dots,L\}} \hat{a}_{v,u}^{i,r}, \quad (12)$$

$$\hat{a}_u^r = \frac{1}{m} \sum_{i \in \{0,1,\dots,L\}} \hat{a}_{u,v}^{i,r},$$

where m indicates the number of stages used. The intermediate mask is defined as follows:

$$M_{inter} = \begin{cases} True, & \text{if } \hat{a}^r > \mu \\ False, & \text{otherwise} \end{cases} \quad (13)$$

where μ is a default value of 0.57, and \hat{a}^r denotes the anomaly score of an image. Then apply collaborative voting with the intermediate mask M_{inter} to refine M .

4. Experiments

4.1. Experimental Setups

Datasets. Our experiments are primarily conducted on the MVTEC AD [4], VisA[48] and BTAD[32] datasets. The MVTEC dataset[4] comprises 15 categories and a total of 5,354 images, of which 3,629 are for training, and 1,725 are for testing, with resolutions varying between 700×700 and $1,024 \times 1,024$ pixels. The VisA dataset[48] provides 12 categories with a total of 10,821 images, including 9,621 normal images and 1,200 anomalous images, with a resolution of $1,500 \times 1,000$ pixels.

Evaluation Metrics. Following previous works[9, 16, 19, 26], we employ the Area Under the Receiver Operating Characteristic Curve (AU-ROC), F1-score at optimal threshold (F_1 -max)[48] and Average Precision (AP)[42] for both anomaly classification and anomaly segmentation. Per-Region Overlap (AU-PRO)[5] is also utilized for anomaly segmentation.

Implementation Details. All test images have been resized to a resolution of 518×518 . We employ the pre-trained CLIP model with ViT-L/14-336[1, 34] as our back-

bone, keeping all parameters frozen throughout the experiments. The model consists of 24 layers, organized into 4 stages with 6 layers each. Patch tokens are extracted from the outputs of 6-th, 12-th, 18-th, and 24-th layers. The CLS token used at the image level is derived from the original CLIP output. In addition, we only alter the final layer of CLIP in Sec.3.2. All experiments are conducted on a single RTX 3090 GPU.

4.2. Main Results

We compare our method with SoTA methods on anomaly classification and segmentation metrics for zero-shot and 1-shot anomaly detection.

Quantitative Results. According to Tab.1, on MVTEC AD dataset[4] we demonstrate a substantial enhancement compared to the training-free WinCLIP, which uses manually designed text prompts. Our method achieved an 8.2% \uparrow and 3.5% \uparrow increase on AU-ROC for anomaly segmentation and classification, respectively, and over a 20% \uparrow improvement in AP and AU-PRO for anomaly segmentation. Additionally, our method outperforms our baseline, MuSc[26], with an 0.7% \uparrow increase in AU-ROC for anomaly segmentation and a 5.6% \uparrow increase for anomaly classification.

Table 1. Comparison with some state-of-the-art on zero-shot and 1-shot methods. Bold values indicate the best results, while underlined values represent the suboptimal results. MuSc-2 denotes inference with 2 images.

Datasets	Methods	Train	Setting	Segmentation				Classification			
				AU-ROC	F_1 -max	AP	AU-PRO	AU-ROC	F_1 -max	AP	
MVTEC	WinCLIP[19]	✗	0-shot	85.1	31.7	18.2	64.6	<u>91.8</u>	92.9	<u>96.5</u>	
	APRIL-GAN[9]	✓	0-shot	87.6	43.3	40.8	44.0	86.1	90.4	93.5	
	AnomalyCLIP[47]	✓	0-shot	91.1	-	-	81.4	91.5	-	96.2	
	AdaCLIP[7]	✓	0-shot	88.7	43.4	-	-	89.2	90.6	-	
	ACR[23]	✓	0-shot	92.5	-	-	-	85.8	-	-	
	[46]	✓	0-shot	92.8	42.5	-	84.0	93.2	94.1	96.7	
	MuSc-2[26]	✗	0-shot	<u>92.6</u>	<u>46.8</u>	<u>42.8</u>	<u>86.4</u>	89.7	<u>93.0</u>	95.2	
	Ours	✗	0-shot	93.3	49.1	46.5	89.6	95.3	95.4	98.3	
	PatchCore[35]	✗	1-shot	93.3	53.0	-	82.3	86.3	92.0	93.8	
	WinCLIP[19]	✗	1-shot	95.2	55.9	-	87.1	<u>93.1</u>	<u>93.7</u>	<u>96.5</u>	
	APRIL-GAN[9]	✓	1-shot	<u>95.1</u>	<u>54.2</u>	51.8	90.6	92.0	92.4	95.8	
	Ours	✗	0-shot	93.3	49.1	46.5	<u>89.6</u>	95.3	95.4	98.3	
	VisA	WinCLIP[19]	✗	0-shot	79.6	14.8	5.4	56.8	78.1	79.0	81.2
		APRIL-GAN[9]	✓	0-shot	94.2	32.3	25.7	86.8	78.0	78.7	81.4
AnomalyCLIP[47]		✓	0-shot	95.5	-	-	<u>87.0</u>	82.1	-	<u>85.4</u>	
AdaCLIP		✓	0-shot	95.5	37.7	-	-	85.8	<u>83.1</u>	-	
[46]		✓	0-shot	94.2	24.1	-	79.7	82.9	80.9	84.7	
MuSc-2[26]		✗	0-shot	<u>95.6</u>	32.6	<u>26.1</u>	80.5	74.0	77.5	78.1	
Ours		✗	0-shot	97.4	<u>35.8</u>	30.8	87.6	<u>85.5</u>	83.5	88.2	
PatchCore[35]		✗	1-shot	95.4	38.0	-	80.5	79.9	81.7	82.8	
WinCLIP[19]		✗	1-shot	<u>96.4</u>	41.3	-	85.1	83.8	83.1	85.1	
APRIL-GAN[9]		✓	1-shot	96.0	<u>38.5</u>	30.9	90.0	91.2	86.9	93.3	
Ours		✗	0-shot	97.4	35.8	<u>30.8</u>	<u>87.6</u>	<u>85.5</u>	<u>83.5</u>	<u>88.2</u>	

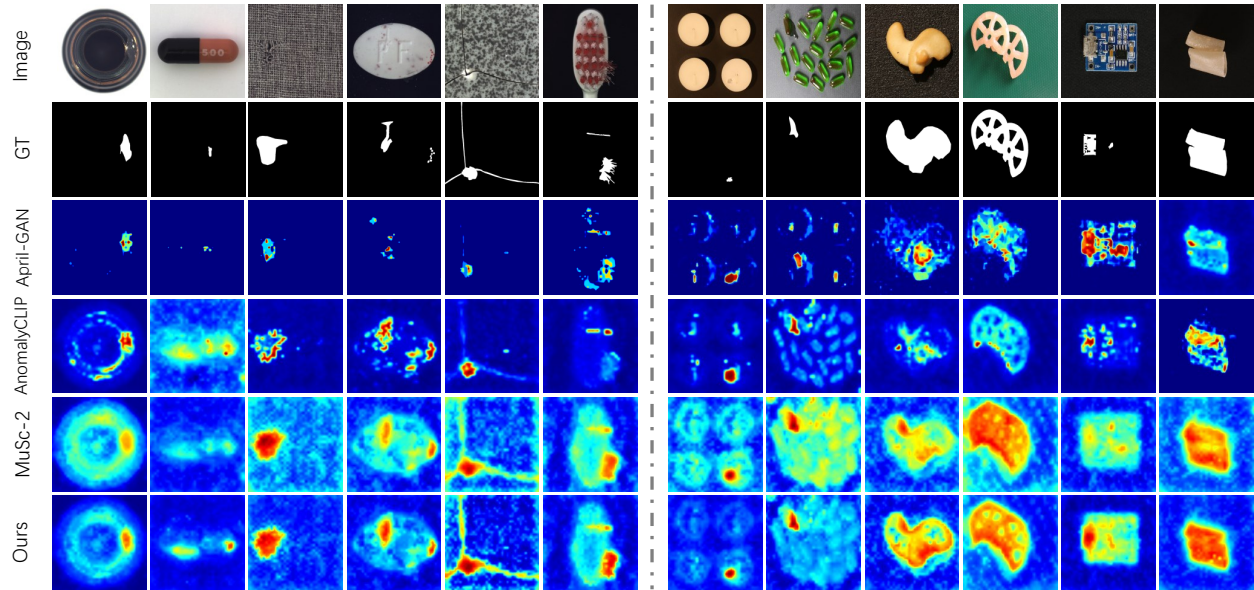


Figure 5. Visualization results of anomaly localization for each class on MVTec AD (left) dataset and VisA dataset (right).

Furthermore, compared to 1-shot methods (using one normal image as a reference), our approach still demonstrates superior performance in anomaly classification. For instance, we achieve improvements of 2.2% \uparrow , 1.7% \uparrow , and 1.8% \uparrow on AU-ROC, F_1 -max, and AP for anomaly classification compared to WinCLIP[19]. On VisA dataset[48], we continue to achieve SOTA performance, especially in anomaly segmentation, where we significantly outperform previous methods and even exceed the segmentation performance of 1-shot anomaly segmentation approaches.

Table 2. Experiments on different noise filtering methods. Bold values indicate the best results

Method	Segmentation				Classification		
	AU-ROC	F_1 -max	AP	AU-PRO	AU-ROC	F_1 -max	AP
w/o filtering	92.8	46.6	42.7	86.5	89.4	92.7	94.9
M	93.1	47.9	44.9	88.2	92.8	93.3	96.8
M^{r-i}	92.8	46.7	42.7	86.6	89.5	92.8	94.9
$M^{i,r}$	93.1	48.1	45.2	87.8	91.7	93.2	96.9
M^{r*i}	93.1	48.2	45.4	88.2	93.0	93.8	97.0
$M^{r,i}$	93.2	48.3	45.5	88.5	93.5	94.7	97.3

Qualitative Results. As shown in the Fig.5, we visualized anomaly segmentation on a subset of data from the MVTec AD[4] and VisA[48] datasets. It can be observed that our method accurately localizes anomalies in most objects, demonstrating the effectiveness of our approach. Compared with the previous methods [9, 26, 47], our method has better anomaly segmentation results.

4.3. Ablation Study

Unless otherwise specified, all ablation experiments were conducted on the MVTec dataset.

Effectiveness comparison of noise filtering and different filtering methods. As shown in Fig.6, we analyze the anomaly classification scores on the MVTec AD dataset[4] with and without noisy feature filtering, with 'Avg' denoting the average anomaly score before and after filtering. The figure shows a clear increase in anomaly scores following filtering, and the quantitative rise in the average score highlights the role of filtering in improving anomaly detection performance.

Additionally, we conduct experiments on different filtering methods. 'W/o filtering' is used to indicate that no noise filtering is applied. As shown in Alg.1, this study implements a two-layer loop structure, where the outer loop iterates over neighborhood F^r , $r \in [1, 2, 3]$, and the inner loop iterates over various layers F^i , $i \in [6, 12, 18, 24]$ for each F^r , denoted as $M^{r,i}$. The outer loop's neighbors, when set in reverse order, are denoted as M^{r-i} . Switching the inner and outer loops is represented by $M^{i,r}$. The notation M^{i*r} indicates a single-layer loop.

Effectiveness comparison of different hyper-parameter.

There are two hyperparameters in Eq.9 and Eq.13 that need to be tuned. This ablation experiment exclusively applies feature matching with noisy feature filtering from Sec.3.3, without incorporating anomaly score adjustments from image-text alignment. As shown in Tab.3 Tab.4, we conduct experiments with different values of λ and μ . The results indicate that the best performance is achieved when λ is and μ is 2. Moreover, the parameters show only a minor

Table 3. Comparative study of the anomaly classification and segmentation performance on MVTEC AD with various λ settings. Bold values indicate the best results.

λ	Segmentation				Classification		
	AU-ROC	F_1 -max	AP	AU-PRO	AU-ROC	F_1 -max	AP
0.95	93.0	47.9	44.8	88.3	93.6	94.6	97.3
1.00	93.1	48.0	44.9	88.3	93.5	94.6	97.3
1.05	93.1	48.0	44.9	88.3	93.5	94.6	97.3
1.10	93.2	48.1	45.0	88.5	93.5	94.7	97.3
1.20	93.1	48.1	44.9	88.3	93.4	94.3	97.3

effect on the experimental outcomes.

Table 5. Experiments on different self-attention mechanisms. Bold values indicate the best results

Attn	Segmentation				Classification		
	AU-ROC	F_1 -max	AP	AU-PRO	AU-ROC	F_1 -max	AP
3	83.2	28.4	22.8	57.8	88.6	90.8	94.8
6	77.3	22.2	15.5	54.8	89.0	91.7	94.6
9	88.8	34.2	27.2	76.5	90.8	91.7	95.8
15	87.3	36.8	32.4	77.4	90.5	91.1	95.8
18	79.1	25.0	18.0	64.0	88.2	90.5	94.7
21	72.1	20.5	12.6	50.9	87.3	91.3	94.6
24	79.5	25.1	18.2	57.5	87.3	91.2	94.4
$v-v$	84.2	34.3	26.0	77.0	89.9	91.2	95.2
$q-q$	80.8	29.7	22.0	65.6	87.1	91.1	94.2
$k-k$	84.2	34.0	25.8	72.6	89.5	91.1	95.4
12 (ours)	88.6	37.5	31.4	83.0	91.6	92.6	96.0

Effectiveness comparison of different attention mechanisms. We achieve anomaly classification and segmentation by calculating the similarity between visual and text features. In the Tab.5, different numbers denote the use of attention weights from the specified layer in place of the final layer’s weights, with $v-v$, $k-k$ and $q-q$ signifying modifications of the final layer’s $q-k$ attention to $v-v$, $k-k$, and $q-q$ attention, respectively. According to Tab.5, replacing the final layer with the attention weights from the twelfth layer achieves the best results on most metrics. Additionally, as depicted in the Fig.3, intermediate layer attention maps emphasize local semantics, whereas deeper layers focus on global tokens. Shallow layers, however, display a more chaotic structure and given that CLIP classifies using the final layer’s output, shallow layers may experience modality shifts.

5. Conclusion

In this paper, we propose a new framework for zero-shot anomaly detection, named FiSeCLIP. We introduce an approach more aligned with practical industrial settings, using images within a batch as mutual references. Since directly using unlabeled test images as references can introduce noise, we fully leverage CLIP’s multi-modal capabilities to filter out noisy features. We then obtain the anomaly

Table 4. Comparative study of the anomaly classification and segmentation performance on MVTEC AD with various μ settings. Bold values indicate the best results.

μ	Segmentation				Classification		
	AU-ROC	F_1 -max	AP	AU-PRO	AU-ROC	F_1 -max	AP
0.47	93.1	48.0	44.7	88.3	93.3	94.2	97.1
0.50	93.2	48.1	44.8	88.3	93.4	94.3	97.2
0.53	93.2	48.1	44.9	88.3	93.4	94.4	97.2
0.55	93.2	48.1	44.9	88.3	93.3	94.3	97.2
0.57	93.2	48.2	45.0	88.3	93.4	94.2	97.2

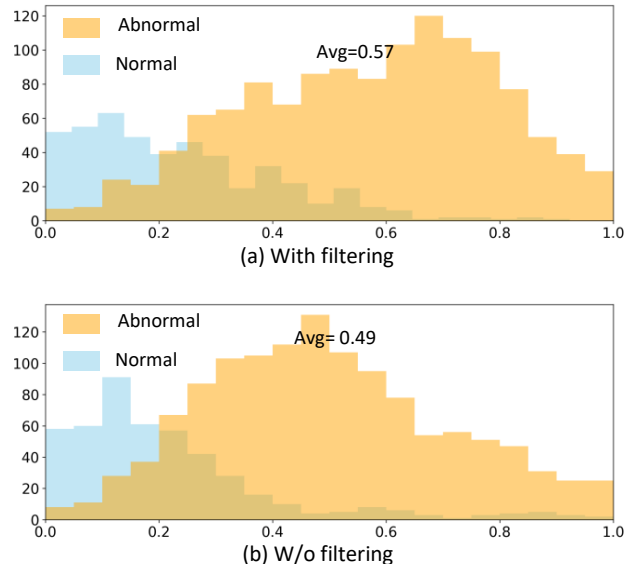


Figure 6. Histogram of anomaly classification scores on the MVTEC AD dataset. (a) represents filtering out noisy features, (b) represents not filtering noisy features, and Avg denotes the average anomaly score.

score through feature matching. Additionally, we aggregate features with varying degrees of neighborhood overlap in a layered calculation, using the anomaly score at each degree to refine the precision of the previous filtering. Furthermore, we enhance CLIP’s self-attention mechanism by substituting the final layer’s attention weights with those of an intermediate layer that better retains local semantics. Comprehensive testing on anomaly detection datasets confirms that our model achieves SoTA performance.

References

- [1] Dosovitskiy Alexey. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 6
- [2] Jaehyeok Bae, Jae-Han Lee, and Seyun Kim. Pni: Industrial anomaly detection using position and neighborhood information. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023. 2
- [3] Yuhu Bai, Jiangning Zhang, Zhaofeng Chen, Yuhang Dong, Yunkang Cao, and Guanzhong Tian. Dual-path frequency discriminators for few-shot anomaly detection. *Knowledge-Based Systems*, 2024. 2
- [4] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. Mvtec ad—a comprehensive real-world dataset for unsupervised anomaly detection. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit*, pages 9592–9600, 2019. 2, 6, 7
- [5] Paul Bergmann, Kilian Batzner, Michael Fauser, David Sattlegger, and Carsten Steger. The mvtec anomaly detection dataset: a comprehensive real-world dataset for unsupervised anomaly detection. *International Journal of Computer Vision*, 2021. 6
- [6] Yunkang Cao, Xiaohao Xu, Jiangning Zhang, Yuqi Cheng, Xiaonan Huang, Guansong Pang, and Weiming Shen. A survey on visual anomaly detection: Challenge, approach, and prospect. *arXiv:2401.16402*, 2024. 1
- [7] Yunkang Cao, Jiangning Zhang, Luca Frittoli, Yuqi Cheng, Weiming Shen, and Giacomo Boracchi. Adaclip: Adapting clip with hybrid learnable prompts for zero-shot anomaly detection. In *European Conference on Computer Vision*, 2025. 2, 3, 6
- [8] Qiyu Chen, Huiyuan Luo, Chengkan Lv, and Zhengtao Zhang. A unified anomaly synthesis strategy with gradient ascent for industrial anomaly detection and localization. *arXiv preprint arXiv:2407.09359*, 2024. 2
- [9] Xuhai Chen, Yue Han, and Jiangning Zhang. April-gan: A zero-/few-shot anomaly classification and segmentation method for cvpr 2023 vand workshop challenge tracks 1&2: 1st place on zero-shot ad and 4th place on few-shot ad. *arXiv preprint arXiv:2305.17382*, 2023. 1, 3, 6, 7
- [10] Xuhai Chen, Jiangning Zhang, Guanzhong Tian, Haoyang He, Wuhao Zhang, Yabiao Wang, Chengjie Wang, and Yong Liu. Clip-ad: A language-guided staged dual-path model for zero-shot anomaly detection. *arXiv preprint arXiv:2311.00453*, 2023. 1, 3
- [11] Thomas Defard, Aleksandr Setkov, Angélique Loesch, and Romaric Audigier. Padim: a patch distribution modeling framework for anomaly detection and localization. In *Proc. Int. Conf. Pattern. Recognit*, 2021. 2
- [12] Hanqiu Deng and Xingyu Li. Anomaly detection via reverse distillation from one-class embedding. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit*, 2022. 1
- [13] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2
- [14] Zhaopeng Gu, Bingke Zhu, Guibo Zhu, Yingying Chen, Hao Li, Ming Tang, and Jinqiao Wang. Filo: Zero-shot anomaly detection by fine-grained description and high-quality localization. *arXiv preprint arXiv:2404.13671*, 2024. 1, 2, 3
- [15] Denis Gudovskiy, Shun Ishizaka, and Kazuki Kozuka. Cflow-ad: Real-time unsupervised anomaly detection with localization via conditional normalizing flows. In *Proc. IEEE Winter Conf. Appl. Comput. Vis.*, 2022. 2
- [16] Haoyang He, Yuhu Bai, Jiangning Zhang, Qingdong He, Hongxu Chen, Zhenye Gan, Chengjie Wang, Xiangtai Li, Guanzhong Tian, and Lei Xie. Mambaad: Exploring state space models for multi-class unsupervised anomaly detection. *arXiv preprint arXiv:2404.06564*, 2024. 2, 6
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 2
- [18] Teng Hu, Jiangning Zhang, Ran Yi, Yuzhen Du, Xu Chen, Liang Liu, Yabiao Wang, and Chengjie Wang. Anomalyd-iffusion: Few-shot anomaly image generation with diffusion model. In *Proc. AAAI Conf. Artif. Intell.*, 2024. 2
- [19] Jongheon Jeong, Yang Zou, Taewan Kim, Dongqing Zhang, Avinash Ravichandran, and Onkar Dabeer. Winclip: Zero-/few-shot anomaly classification and segmentation. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit*, pages 19606–19616, 2023. 1, 3, 6, 7
- [20] Mengcheng Lan, Chaofeng Chen, Yiping Ke, Xinjiang Wang, Litong Feng, and Wayne Zhang. Clearclip: Decomposing clip representations for dense vision-language inference. *arXiv preprint arXiv:2407.12442*, 2024. 2, 3
- [21] Mengcheng Lan, Chaofeng Chen, Yiping Ke, Xinjiang Wang, Litong Feng, and Wayne Zhang. Proxyclick: Proxy attention improves clip for open-vocabulary segmentation. *arXiv preprint arXiv:2408.04883*, 2024. 2
- [22] Jiarui Lei, Xiaobo Hu, Yue Wang, and Dong Liu. Pyramid-flow: High-resolution defect contrastive localization using pyramid normalizing flow. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit*, pages 14143–14152, 2023. 1
- [23] Aodong Li, Chen Qiu, Marius Kloft, Padhraic Smyth, Maja Rudolph, and Stephan Mandt. Zero-shot anomaly detection via batch normalization. In *NeurIPS*, 2023. 6
- [24] Aodong Li, Chen Qiu, Marius Kloft, Padhraic Smyth, Maja Rudolph, and Stephan Mandt. Zero-shot anomaly detection via batch normalization. *Advances in Neural Information Processing Systems*, 2023. 1
- [25] Chun-Liang Li, Kihyuk Sohn, Jinsung Yoon, and Tomas Pfister. Cutpaste: Self-supervised learning for anomaly detection and localization. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit*, 2021. 2
- [26] Xurui Li, Ziming Huang, Feng Xue, and Yu Zhou. Musc: Zero-shot industrial anomaly classification and segmentation with mutual scoring of the unlabeled images. *arXiv preprint arXiv:2401.16753*, 2024. 1, 2, 3, 5, 6, 7
- [27] Yi Li, Hualiang Wang, Yiqun Duan, and Xiaomeng Li. Clip surgery for better explainability with enhancement in open-vocabulary tasks. *arXiv preprint arXiv:2304.05653*, 2023. 2
- [28] Yufei Liang, Jiangning Zhang, Shiwei Zhao, Runze Wu, Yong Liu, and Shuwen Pan. Omni-frequency channel-

- selection representations for unsupervised anomaly detection. *IEEE Trans. Image Process.*, 2023. [2](#)
- [29] 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*, 2024. [1](#)
- [30] Jiaqi Liu, Guoyang Xie, Jinbao Wang, Shangnian Li, Chengjie Wang, Feng Zheng, and Yaochu Jin. Deep industrial image anomaly detection: A survey. *Mach. Intell. Res.*, 21:104–135, 2024. [1](#)
- [31] Shuai Lyu, Dongmei Mo, and Wai keung Wong. Reb: Reducing biases in representation for industrial anomaly detection. *Knowledge-Based Syst.*, 290:111563, 2024. [2](#)
- [32] Pankaj Mishra, Riccardo Verk, Daniele Fornasier, Claudio Piciarelli, and Gian Luca Foresti. Vt-adl: A vision transformer network for image anomaly detection and localization. In *2021 IEEE 30th International Symposium on Industrial Electronics (ISIE)*, 2021. [6](#)
- [33] Zhen Qu, Xian Tao, Mukesh Prasad, Fei Shen, Zhengtao Zhang, Xinyi Gong, and Guiguang Ding. Vcp-clip: A visual context prompting model for zero-shot anomaly segmentation. *arXiv preprint arXiv:2407.12276*, 2024. [3](#)
- [34] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. [1](#), [2](#), [6](#)
- [35] Karsten Roth, Latha Pemula, Joaquin Zepeda, Bernhard Schölkopf, Thomas Brox, and Peter Gehler. Towards total recall in industrial anomaly detection. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit*, 2022. [1](#), [2](#), [5](#), [6](#)
- [36] João Santos, Triet Tran, and Oliver Rippel. Optimizing patchcore for few/many-shot anomaly detection. *arXiv:2307.10792*, 2023. [5](#)
- [37] Tong Shao, Zhuotao Tian, Hang Zhao, and Jingyong Su. Explore the potential of clip for training-free open vocabulary semantic segmentation. In *European Conference on Computer Vision*, pages 139–156. Springer, 2025. [2](#)
- [38] Shashanka Venkataramanan, Kuan-Chuan Peng, Rajat Vikram Singh, and Abhijit Mahalanobis. Attention guided anomaly localization in images. In *European Conference on Computer Vision*, 2020. [1](#), [2](#)
- [39] Feng Wang, Jieru Mei, and Alan Yuille. Sclip: Rethinking self-attention for dense vision-language inference. In *European Conference on Computer Vision*, 2025. [2](#)
- [40] Julian Wyatt, Adam Leach, Sebastian M Schmon, and Chris G Willcocks. Anoddpn: Anomaly detection with denoising diffusion probabilistic models using simplex noise. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit*, pages 650–656, 2022. [2](#)
- [41] Zhiyuan You, Lei Cui, Yujun Shen, Kai Yang, Xin Lu, Yu Zheng, and Xinyi Le. A unified model for multi-class anomaly detection. In *Proc. Adv. Neural Inf. Process. Syst.*, pages 4571–4584, 2022. [1](#)
- [42] Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. Draema: a discriminatively trained reconstruction embedding for surface anomaly detection. In *Proc. IEEE Int. Conf. Comput. Vis.*, 2021. [2](#), [6](#)
- [43] Jiangning Zhang, Xuhai Chen, Yabiao Wang, Chengjie Wang, Yong Liu, Xiangtai Li, Ming-Hsuan Yang, and Dacheng Tao. Exploring plain vit reconstruction for multi-class unsupervised anomaly detection. *arXiv:2312.07495*, 2023. [2](#)
- [44] Jiangning Zhang, Haoyang He, Zhenye Gan, Qingdong He, Yuxuan Cai, Zhucun Xue, Yabiao Wang, Chengjie Wang, Lei Xie, and Yong Liu. A comprehensive library for benchmarking multi-class visual anomaly detection. *arXiv preprint arXiv:2406.03262*, 2024. [1](#)
- [45] Ximiao Zhang, Min Xu, and Xiuzhuang Zhou. Realnet: A feature selection network with realistic synthetic anomaly for anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. [2](#)
- [46] Zhaoxiang Zhang, Hanqiu Deng, Jinan Bao, and Xingyu Li. Dual-image enhanced clip for zero-shot anomaly detection. *arXiv preprint arXiv:2405.04782*, 2024. [3](#), [6](#)
- [47] Qihang Zhou, Guansong Pang, Yu Tian, Shibo He, and Jiming Chen. Anomalyclip: Object-agnostic prompt learning for zero-shot anomaly detection. *arXiv preprint arXiv:2310.18961*, 2023. [1](#), [2](#), [3](#), [6](#), [7](#)
- [48] Yang Zou, Jongheon Jeong, Latha Pemula, Dongqing Zhang, and Onkar Dabeer. Spot-the-difference self-supervised pre-training for anomaly detection and segmentation. In *Proc. Eur. Conf. Comput. Vis.*, pages 392–408. Springer, 2022. [2](#), [6](#), [7](#)