# AdaCLIP: Adapting CLIP with Hybrid Learnable Prompts for Zero-Shot Anomaly Detection

Yunkang Cao<sup>1,2</sup>, Jiangning Zhang<sup>3,4</sup>, Luca Frittoli<sup>2</sup>, Yuqi Cheng<sup>1</sup>, Weiming Shen<sup>1⊠</sup>, and Giacomo Boracchi<sup>2</sup>

<sup>1</sup> Huazhong University of Science and Technology {cyk\_hust,yuqicheng,shenwm}@hust.edu.cn <sup>2</sup> Politecnico di Milano {luca.frittoli,giacomo.boracchi}@polimi.it <sup>3</sup> Zhejiang University <sup>4</sup> Youtu Lab, Tencent 186368@zju.edu.cn

Abstract. Zero-shot anomaly detection (ZSAD) targets the identification of anomalies within images from arbitrary novel categories. This study introduces AdaCLIP for the ZSAD task, leveraging a pre-trained visionlanguage model (VLM), CLIP. AdaCLIP incorporates learnable prompts into CLIP and optimizes them through training on auxiliary annotated anomaly detection data. Two types of learnable prompts are proposed: static and dynamic. Static prompts are shared across all images, serving to preliminarily adapt CLIP for ZSAD. In contrast, dynamic prompts are generated for each test image, providing CLIP with dynamic adaptation capabilities. The combination of static and dynamic prompts is referred to as hybrid prompts, and yields enhanced ZSAD performance. Extensive experiments conducted across 14 real-world anomaly detection datasets from industrial and medical domains indicate that AdaCLIP outperforms other ZSAD methods and can generalize better to different categories and even domains. Finally, our analysis highlights the importance of diverse auxiliary data and optimized prompts for enhanced generalization capacity. Code is available at https://github.com/caoyunkang/AdaCLIP.

Keywords: Anomaly Detection  $\cdot$  Prompt Learning  $\cdot$  Zero-shot Learning

# 1 Introduction

Anomaly detection (AD) in images [12, 13] holds significant importance across various domains, including industrial inspection [3,33,48] and medical diagnosis [7]. The primary goal of AD methods is to detect deviations from normal patterns, either image or pixel-level. Most AD methods rely on unsupervised learning [9,41] and semi-supervised learning [11, 17, 42] paradigms that require either normal samples or annotated anomalous samples from the target category for training, as depicted in Fig. 1. For instance, to train a dedicated model for the category 'chewing gum', traditional unsupervised AD methods require a substantial dataset

comprising normal 'chewing gum' images, while semi-supervised approaches impose an even stricter requirement, requiring annotated abnormal images.

Some scenarios are characterized by the *cold start* problem, meaning that it is not feasible to gather enough normal images for training an unsupervised model, thus preventing both unsupervised and semi-supervised AD solutions. The emerging zero-shot anomaly detection (ZSAD [24]) paradigm addresses this issue, aiming at detecting anomalies in images belonging to unseen categories, without requiring any image of that category for training. Existing ZSAD methods commonly rely on pre-trained vision-language models (VLMs) due to their broad generalization capability. Some ZSAD methods employ VLMs for ZSAD without any additional training [24, 46], while others leverage annotated images from auxiliary anomaly-detection datasets to tailor VLMs for ZSAD, as Fig. 1 shows.

The pioneering ZSAD method, WinCLIP [24], directly uses pre-trained VLMs with hand-crafted textual prompts to identify anomalies. Similarly to zero-shot classification, WinCLIP detects as anomalous images that are close to the selected prompts in the embedding space. However, WinCLIP exhibits limited detection performance since its underlying VLM, CLIP [40], is trained on natural image-text datasets [43] and is not specialized for anomaly detection. Conversely, APRIL-GAN [14] and AnomalyCLIP [56] address ZSAD by adapting VLMs on auxiliary anomaly-detection datasets that contain annotated anomalies. This adaptation scheme is gaining popularity due to the growing availability of annotated AD datasets [3,57] spanning diverse categories [3] and domains [18,57]. Importantly, the adaptation scheme adheres to the zero-shot learning paradigm, as long as testing images do not belong to categories presented in the auxiliary AD dataset.

The rationale behind ZSAD approaches is that testing images may exhibit universal patterns, either normal or anomalous, that VLMs can identify. Additionally, adapting VLMs on auxiliary data can be beneficial as these data might contain patterns that are useful for detecting anomalies in novel categories. For example, the scratches on 'pill' images might improve the model's ability to detect similar abnormal patterns on 'chewing gum' (as illustrated in Fig. 1).

To take the most from auxiliary datasets for ZSAD, we propose AdaCLIP, which builds upon the mainstream zero-shot learning principle in CLIP. In particular, AdaCLIP computes similarities between patch embeddings and text embeddings for textual captions describing normal/abnormal states using CLIP. To enhance the ZSAD performance, AdaCLIP introduces additional lightweight learnable parameters in two forms: *projection* and *prompting* layers. As in APRIL-GAN [14], our projection layer is designed to align the dimensions between patch tokens and text embeddings, while introducing additional learnable parameters for fine-tuning CLIP. Prompting layers are used to replace the original transformer layers within CLIP, by concatenating additional prompting tokens and the layer input. Prompting has proven very effective in adapting VLMs [29]. To ease the adaptation with auxiliary data, static and dynamic learnable prompts are introduced, where static prompts are shared across all images and dynamic prompts are generated based on the testing image. The combination of static



Fig. 1: Left: Illustrations for training and test data of unsupervised, semi-supervised, and zero-shot anomaly detection paradigms. Right: Quantitative comparison with popular methods by pixel-level max-F1 [24] on industrial and medical datasets.

and dynamic prompts, referred to as hybrid prompts, demonstrates significant generalization capabilities and promising ZSAD performance, as shown in Fig. 1. In summary, our contributions include the following key components:

- We introduce a novel ZSAD method named AdaCLIP. AdaCLIP comprises hybrid (static and dynamic) learnable prompts to better exploit the auxiliary data to enhance ZSAD performance. A hybrid-semantic fusion module is also developed to extract region-level context about anomaly regions, thereby enhancing image-level anomaly detection performance.
- We show that different VLMs –not only CLIP– can be effectively adapted for ZSAD. Additionally, we demonstrate the importance of optimized prompts for detecting anomalies within individual images.

Our experiments demonstrate that we achieve state-of-the-art (SOTA) performance in ZSAD across 14 datasets spanning industrial and medical domains. We showcase that our AdaCLIP can effectively leverage information from auxiliary datasets, even when referring to categories from different domains (medical/industrial), outperforming alternative ZSAD methods. Additionally, we underscore that leveraging diverse auxiliary data is beneficial for ZSAD.

# 2 Related Work

#### 2.1 Traditional Anomaly Detection

**Unsupervised Anomaly Detection** methods like [8,41] learn exclusively from normal samples within target categories. Unsupervised AD methods typically model normal sample distributions during training and subsequently compare test samples to the learned normal sample distribution to detect anomalies. A very effective approach consists in extracting features from each sample using pre-trained neural networks [6,9,16], and then modeling the features distribution by knowledge distillation [27, 34, 53], reconstruction [5, 22, 50, 51], or memory bank-based approaches [23, 47].

**Semi-supervised Anomaly Detection** methods like [11, 17] require both normal and abnormal images with annotations from target categories for training. They typically utilize annotated abnormal samples to learn a more compact description boundary for normal samples. Since some additional abnormal samples are exploited, they typically present better AD performance in comparison to unsupervised AD but impose a strict requirement for data.

Despite the promising anomaly detection performance achieved by these traditional AD methods, their effectiveness tends to diminish when fewer normal samples are available for training. In contrast, we aim to develop a generic ZSAD model for anomaly detection across unseen categories without training samlpes.

#### 2.2 Zero-shot Anomaly Detection

Zero-shot learning often requires extensive training data to attain generalization abilities [15,24]. Many off-the-shelf VLMs have been developed, presenting promising zero-shot capabilities. These pre-trained VLMs are leveraged to identify anomalies across unbounded categories. For instance, WinCLIP [24] employs CLIP [40] to compute similarities between embeddings of image patches and embeddings of captions regarding normal/abnormal states, which is subsequently enhanced by text augmentation in [46]. In contrast, SAA [10] utilizes Grounding DINO [35] to identify abnormal regions within a test image using text prompts, followed by refinement with SAM [30]. However, these VLMs are typically trained on natural image-text pairs and are not specifically designed for AD. Therefore, APRIL-GAN [14] and CLIP-AD [15] enhance the ZSAD performance of CLIP by tuning additional projection layers with annotated auxiliary AD data. With these auxiliary data, AnomalyCLIP [56] preliminarily explores prompt learning and introduces learnable text prompts to adapt VLMs for ZSAD. AnomalyGPT [19] also introduces textual prompting learning but for unsupervised AD. In this paper, we further delve into prompt learning and develop multimodal hybrid learnable prompts to maximize the utility of auxiliary AD data. Table 4 in Appendix highlights the significance of AdaCLIP in comparison to other alternatives.

#### 2.3 Prompt Learning

In the realm of VLMs, prompt learning [29] involves incorporating learnable tokens into the input image or text, effectively tailoring VLMs to specific scenarios. Early prompt learning methods introduce static prompts to VLMs. For instance, CoOp [55] integrates learnable tokens in addition to the input text into the text branch. However, recent advancements in prompt learning methods [52] have identified that static prompts may be susceptible to distribution diversity. Consequently, CoCoOp [54] and IDPT [52] propose generating dynamic prompts based on the inputs for improving modeling capabilities. Whereas previous prompt learning methods primarily focused on the text encoder of VLMs [54],

<sup>4</sup> Y. Cao et al.



Fig. 2: Framework of AdaCLIP.

recent studies [26,29] have increasingly acknowledged the significance of prompting the image encoder, *i.e.*, visual prompting, to better exploit the multimodal capabilities of VLMs. In this paper, we propose multimodal (image+text) hybrid (static+dynamic) prompts to adapt VLMs for improving anomaly detection.

## 3 Problem Formulation

Our objective is to develop a model that associates an input image  $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$ with an image-level anomaly score S and a pixel-level anomaly map  $\mathbf{M} \in \mathbb{R}^{H \times W}$ , indicating whether  $\mathbf{I}$  and its pixels are normal or abnormal. Typically, the values of the anomaly score and anomaly map pixels fall within the range [0, 1], where larger values indicate higher probabilities of being abnormal. We operate within the ZSAD context, training our model using an auxiliary anomaly detection dataset  $\mathcal{I}_{\text{train}} = \{(\mathbf{I}_i, \mathbf{G}_i)\}_{i=1}^{N_{\text{train}}}$ , which contains categories distinct from those in the testing dataset  $\mathcal{I}_{\text{test}} = \{\mathbf{I}_i\}_{i=1}^{N_{\text{test}}}$ . The auxiliary training dataset includes both normal and abnormal images  $\mathbf{I}$  along with their annotated masks  $\mathbf{G} \in \mathbb{R}^{H \times W}$ , where pixels have value 0 if normal and 1 if abnormal. By learning from this auxiliary dataset  $\mathcal{I}_{\text{train}}$ , the model is expected to learn normal and abnormal patterns that are common to different classes, enabling the detection of anomalies in novel categories.

## 4 AdaCLIP

#### 4.1 Overview

The framework of AdaCLIP is illustrated in Fig. 2. Given an image I, Ada-CLIP follows the general ZSAD principle of comparing CLIP embedding as WinCLIP [24] do. In particular, we detect anomalies by calculating similarities in CLIP embedding space between the image and textual captions for normal/abnormal states, such as "A photo of normal [CLS]" and "A photo of

damaged [CLS]", where [CLS] denotes to the name of the testing category, like 'carpet', 'hazelnut', etc. Notably, AdaCLIP enhances the pre-trained CLIP by incorporating learnable parameters through prompting layers for image and text encoders, denoted as  $L_I^P$  and  $L_T^P$  respectively, which replace the original transformer layers. For the prompting layers, both static prompts  $\mathbf{P}^S$  and dynamic prompts  $\mathbf{P}^D$  are introduced. AdaCLIP also introduces a projection layer Proj at the end of the image encoder, and a Hybrid Semantic Fusion (HSF) module designed to extract semantic-rich image embeddings for computing image-level anomaly scores S.

## 4.2 Prompting Layers

AdaCLIP introduces prompting layers  $L_I^P$  and  $L_T^P$  to replace the original transformer layers in the image and text encoders of CLIP, respectively. Prompting layers [29] preserves the weights of the transformer (to inherit its generalization ability) but concatenates learnable prompting tokens **P** to the vanilla tokens derived from the input images or texts, as illustrated in Fig. 2. Thanks to the self-attention mechanism in transformer layers, the learnable prompting token will contribute to all the output tokens, including the vanilla ones.

More specifically, prompting tokens  $\mathbf{P} \in \mathbb{R}^{M \times C}$  are concatenated to the input vanilla tokens  $\mathbf{T} \in \mathbb{R}^{N \times C}$  of the transformer layer. Here, C denotes the embedding dimension, while N and M denote the lengths of vanilla tokens and prompting tokens, respectively, where  $M \ll N$  for lightweight adaptation. Let  $L_i^P$  denote the *j*-th prompting layer, then the feed-forward process is,

$$[\mathbf{T}_{j+1}, \_] = L_j^P([\mathbf{T}_j, \mathbf{P}_j]), \quad j \le J, \tag{1}$$

$$[\mathbf{T}_{j+1}, \mathbf{P}_{j+1}] = L_j^P([\mathbf{T}_j, \mathbf{P}_j]), \quad j > J,$$

$$(2)$$

where  $[\cdot, \cdot]$  denotes concatenation along rows. Learnable prompting tokens are incorporated up to a limited depth J, while prompting tokens for the remaining layers are generated through feed-forwarding. Typically, J is set to a small value, as too many learnable parameters may result in overfitting on auxiliary data.

#### 4.3 Hybrid Learnable Prompts

To effectively utilize auxiliary data for enhanced anomaly detection performance, we introduce both *static* and *dynamic* prompts.

Static Prompts  $\mathbf{P}^{S}$ . Static prompts  $\mathbf{P}^{S}$  serve as foundational learning tokens shared across all images, which are explicitly learned from auxiliary data during training, as Fig 2 shows. However, their limited adaptation effectiveness is acknowledged by previous studies [52].

**Dynamic Prompts P**<sup>D</sup>. We further introduce dynamic prompts **P**<sup>D</sup> to enhance the modeling capacity for diverse distributions. Dynamic prompts differ from static prompts as they are generated on each testing image by the Dynamic Prompt Generator (DPG):

$$\mathbf{P}^D = \mathrm{DPG}(\mathbf{I}). \tag{3}$$

In our case DPG is a frozen pre-trained backbone such as CLIP to extract class tokens, followed by a learnable linear layer to project the class tokens into dynamic prompts  $\mathbf{P}^{D}$ . Both dynamic prompts for  $L_{I}^{P}$  in the image encoder and  $L_{T}^{P}$  in the text encoder are generated from the testing image, as shown in Fig. 2.

AdaCLIP sums up the static and dynamic prompts, referred to as hybrid prompts, for both prompting layers  $L_I^P$  and  $L_T^P$ . By replacing the original transformer layers with these prompting layers, the image encoder extracts patch embeddings  $\mathbf{F}^P = \{\mathbf{F}_0^P, \ldots\}$  for the input image I from multiple prompting layers, while the text encoder generates normal/abnormal text embeddings  $\mathbf{F}_N^T$ ,  $\mathbf{F}_A^T$  for the corresponding textual captions.

#### 4.4 Projection Layer

The original CLIP [40] architecture makes the dimensions of patch embeddings and text embeddings unmatched, thus we append a projection layer Proj to the image encoder. In particular, we align the dimensions between patch embeddings ( $\mathbf{F}^{P}$ ) and the embeddings of normal ( $\mathbf{F}_{N}^{T}$ ) and anomalous ( $\mathbf{F}_{A}^{T}$ ) texts by introducing a linear layer with bias. In addition, the projection layer adds some learnable parameters for CLIP adaption.

#### 4.5 Pixel-Level Anomaly Localization

We derive the anomaly score by measuring the cosine similarities between patch embeddings  $\mathbf{F}^{P}$ , and text embeddings  $\mathbf{F}_{N}^{T}$  and  $\mathbf{F}_{A}^{T}$ . We adopt the same approach as in WinCLIP [24], and define the anomaly map from *i*-th layer as follows:

$$\mathbf{M}_{i} = \phi \left( \frac{\exp(\cos(\mathbf{F}_{i}^{P}, \mathbf{F}_{A}^{T}))}{\exp(\cos(\mathbf{F}_{i}^{P}, \mathbf{F}_{N}^{T})) + \exp(\cos(\mathbf{F}_{i}^{P}, \mathbf{F}_{A}^{T}))} \right),$$
(4)

where  $\cos(\cdot, \cdot)$  denotes the cosine similarity and  $\phi$  is a reshape and interpolate function, transforming anomaly scores for patch embeddings into anomaly maps  $\mathbf{M}_i \in \mathbb{R}^{H \times W}$ . Then we take anomaly maps from several layers in a multi-hierarchy manner [24] and aggregate these anomaly maps into a final prediction  $\mathbf{M}$ . During training, AdaCLIP optimizes the pixel-level anomaly map  $\mathbf{M}$  with dice loss [37] and focal loss [32] on the auxiliary data.

#### 4.6 Hybrid Semantic Fusion Module

AdaCLIP introduces an HSF module to improve image-level AD performance. Traditional AD methods [9,14] for image-level AD often select the maximum values of anomaly maps as anomaly scores, but this is sensitive to noisy predictions. In contrast, we present the HSF module to aggregate patch embeddings that are more likely to represent abnormalities, thereby aggregating region-level information for robust image-level anomaly detection. We refer to the output of HSF as semantic-rich embedding  $\mathbf{F}^{I}$ .

As Fig. 2 shows, the HSF module follows a three-step paradigm: ① Cluster patch embeddings into K groups using KMeans [2]. ② Compute the anomaly scores of individual clusters by averaging the scores of the corresponding positions in the anomaly map **M**. ③ Select the cluster with the highest anomaly scores, calculate its centroids, and aggregate them into the final semantic-rich image embedding  $\mathbf{F}^{I}$ , which encapsulates semantic information about the most abnormal clusters. The resulting semantic-rich image embedding effectively improves imagelevel AD performance compared to the maximum value-based anomaly detection methods. More details regarding HSF are presented in Appendix Section 2.2.

#### 4.7 Image-Level Anomaly Detection

After extracting the semantic-rich image embeddings  $\mathbf{F}^{I}$ , we compute the imagelevel anomaly scores S similar to (4), using cosine similarities between  $\mathbf{F}^{I}$  and the text embeddings  $\mathbf{F}_{A}^{N}$  and  $\mathbf{F}_{A}^{T}$ , followed by softmax normalization. Then we optimize image-level anomaly scores S using focal loss [32].

## 5 Experiments

#### 5.1 Experimental Setup

**Datasets.** We conduct experiments using datasets from industrial and medical domains. Specifically, for the industrial domain, we use MVTec AD [3], VisA [57], MPDD [25], BTAD [38], KSDD [44], DAGM [49], and DTD-Synthetic [1] datasets. In the medical domain, we consider brain tumor detection datasets HeadCT [31], BrainMRI [28], Br35H [21], skin cancer detection dataset ISIC [20], colon polyp detection datasets ClinicDB [4], and ColonDB [45], as well as thyroid nodule detection dataset TN3K [18]. A detailed introduction to these datasets can be found in Appendix Section 1.

**Evaluation Metrics.** Following previous ZSAD studies [14, 24], we employ the Area Under the Receiver Operating Characteristic Curve (AUROC) and the maximum F1 score (max-F1) under the optimal threshold to evaluate both image-level and pixel-level AD performance. In addition to dataset-level results, we also report domain-level average performance in the form of (AUROC, max-F1).

**Implementation Details.** This study employs the pre-trained CLIP (ViT-L/14@336px)<sup>5</sup> as the default backbone and extracts patch embeddings from the 6-th, 12-th, 18-th, and 24-th layers. All images undergo resizing to a resolution of  $518 \times 518$  for both training and testing. For the ZSAD task, it is imperative that the auxiliary data does not contain any categories present in the test set. Although ClinicDB [4] and ColonDB [45] both comprise colon polyp data, their appearances differ significantly. Therefore, we default to using the industrial dataset, MVTec AD [3], and the medical dataset, ClinicDB [4], as auxiliary data. For evaluations on MVTec AD and ClinicDB, VisA [57] and ColonDB [45] are utilized for training. The prompting depth J is set to four and the prompting length M is set to five

<sup>&</sup>lt;sup>5</sup> https://github.com/mlfoundations/open clip

by default. We train AdaCLIP for five epochs with a learning rate of 0.01. All experiments are conducted using PyTorch-1.9.2 with a single NVIDIA A6000 48GB GPU. Appendix Section 3 presents further implementation details.

#### 5.2 Main Experimental Results

**Comparison Methods.** This study compares the proposed AdaCLIP with two sets of methods: with and without training on auxiliary data. For methods without training, we reproduce SAA [10] and WinCLIP [24] for comparisons. Regarding methods with training, we choose the existing ZSAD method based on CLIP, APRIL-GAN [14], and AnomalyCLIP [56]. In addition, to explore whether other VLMs excluding CLIP can be adapted for ZSAD, we train DINOV2 [39] and SAM [30] on the auxiliary data by adding linear layers as segmentation heads after multiple transformer layers. More details about the implementation of these methods can be found in Appendix Section 3. Unfortunately, we cannot directly compare with AnomalyCLIP [56] because its implementation is not publicly available before our submission date. To enable a fair comparison, we have evaluated AdaCLIP under the experimental setting of AnomalyCLIP, and the results are reported in Appendix Section 4.

**Zero-shot Anomaly Detection in the Industrial Domain:** Table 1 reports the results in the industrial domain. It distinctly illustrates that methods with training exhibit superior performance compared to alternative ZSAD methods without training on auxiliary data. In particular, WinCLIP and SAA which utilize hand-crafted textual prompts present subpar AD performance. Conversely, adapting DINOV2 and SAM with auxiliary data demonstrates promising pixellevel ZSAD performance. The superior performance of the set of ZSAD methods trained with the auxiliary data underscores that pre-trained VLMs are already endowed with essential knowledge for anomaly detection. This existing knowledge can be effectively leveraged for ZSAD through proper adaptation, like the strategy we employed.

Moreover, as evident in Table 1, the proposed AdaCLIP showcases significant improvements over other ZSAD methods, *e.g.*, 3.7% image-level and 3.3% pixellevel enhancements on max-F1 compared to the second-place method. Also, AdaCLIP achieves the best overall ranking across all datasets in terms of both image- and pixel-level performance. This showcases the excellence of AdaCLIP and validates the efficacy of the introduced prompting layers. We further present visualizations of the predicted anomaly maps across various datasets in Fig. 3. AdaCLIP exhibits significantly more accurate segmentation for novel industrial categories in comparison to other methods. The precise detection results for challenging categories such as tubes, capsules, and pipe fryum further highlight the superiority of AdaCLIP.

**Zero-shot Anomaly Detection in the Medical Domain.** We also conduct experiments in the medical domain to further investigate the generalization ability of these ZSAD methods. The results exhibit a similar trend to those in the industrial domain, where methods with training outperform SAA and WinCLIP by a significant margin. AdaCLIP emerges as the top performer with

Table 1: Comparisons of ZSAD methods in the industrial domain. The best performance is in **bold**, and the second-best is <u>underlined</u>. <sup>†</sup> denotes to results taken from original papers. Rank denotes to the average performance rankings of different methods on various datasets.

Metric	Dataset	w/o supervised training		w/ supervised training			
1100110	Durabet	SAA [10]	WinCLIP [24]	DINOV2 [39]	SAM [30]	APRIL-GAN [14]	AdaCLIP
	MVTec AD	(63.5, 87.4)	$(91.8,92.9)^\dagger$	(74.4, 87.4)	(70.8, 86.0)	(82.3, 88.9)	$(\underline{89.2}, \underline{90.6})$
	VisA	(67.1, 75.9)	$(78.1, 80.7)^{\dagger}$	(75.2, 78.5)	(61.9, 73.9)	$(\underline{81.7}, 80.7)$	(85.8, 83.1)
	MPDD	(42.7, 73.9)	$(61.4, \underline{77.5})$	(62.4, 74.9)	(63.0, 77.0)	$(\underline{66.0}, 76.0)$	(76.0, 82.5)
Image-level	BTAD	(59.0, 89.7)	(68.2, 67.6)	(79.3, 69.3)	(89.4, 85.7)	(85.2, 82.0)	(88.6, 88.2)
(AUROC, max-F1)	KSDD	(68.6, 37.6)	(93.3, 79.0)	(94.9, 77.5)	(65.8, 37.9)	(95.7, 85.2)	(97.1, 90.7)
	DAGM	(87.1, 88.8)	(91.7, 87.6)	(90.7, 89.2)	(82.7, 83.6)	(93.5, 91.8)	(99.1, 97.5)
	DTD-Synthetic	(94.4, 93.5)	(95.1, 94.1)	(85.8, 93.5)	(81.9, 91.1)	$(\overline{98.1}, \overline{96.8})$	$(\underline{95.5},  \underline{94.7})$
	Average	(68.9, 78.1)	(82.8, 82.8)	(80.4, 81.5)	(73.6, 76.4)	( <u>86.1</u> , <u>85.9</u> )	(90.2, 89.6)
	Rank	(5.3, 4.4)	(3.4, 3.4)	(4.0, 4.1)	(4.7, 5.0)	$(\underline{2.1},  \underline{2.6})$	(1.4, 1.4)
	MVTec AD	(75.5, 38.1)	$(85.1, 31.6)^{\dagger}$	(85.9, 39.6)	(85.4, 29.4)	(83.7, <u>39.8</u> )	(88.7, 43.4)
	VisA	(76.5, 31.6)	$(79.6, 14.8)^{\dagger}$	(95.0, 30.3)	(92.6, 18.2)	(95.2, 32.3)	(95.5, 37.7)
	MPDD	(81.7, 18.9)	(71.2, 15.4)	(95.6, 31.1)	(94.8, 22.1)	(95.1, 30.6)	(96.1, 34.9)
Pixel-level	BTAD	(65.8, 14.8)	(72.6, 18.5)	(91.9, 43.4)	(93.8, 46.9)	(89.5, 38.4)	(92.1, 51.7)
(AUROC, max-F1)	KSDD	(78.8, 6.6)	(95.8, 21.3)	(99.3, 50.6)	(91.2, 18.4)	(98.2, 56.2)	(97.7, 54.5)
	DAGM	(62.7, 32.6)	(81.3, 13.9)	(90.9, 52.0)	(88.6, 40.7)	(90.3, 57.9)	(91.5, <u>57.5</u> )
	DTD-Synthetic	(76.7, 60.6)	(79.5, 16.1)	$(\overline{97.0}, 63.4)$	(95.0, 56.7)	( <u>97.8</u> , <b>72.7</b> )	( <b>97.9</b> , <u>71.6</u> )
	Average	(73.9, 29.0)	(80.7, 18.8)	$(\underline{93.7}, 44.3)$	(91.7, 33.2)	$(92.8, \underline{46.9})$	(94.2, 50.2)
	Rank	(5.9, 4.7)	(4.9, 5.6)	$(\underline{2.3}, 3.0)$	(3.6, 4.3)	$(3.0, \underline{2.0})$	(1.4, 1.4)

Table 2: Comparisons of ZSAD methods in the medical domain. The best performance is in **bold**, and the second-best is <u>underlined</u>. Rank denotes to the average performance rankings of different methods on various datasets.

Metric	Metric Dataset		w/o supervised training		w/ supervised training			
1100110	Dataset	SAA [10]	WinCLIP [24]	DINOV2 [39]	SAM [30]	APRIL-GAN [14]	AdaCLIP	
Image-level (AUROC, max-F1)	HeadCT BrainMRI Br35H	$\begin{array}{c} (46.8,68.0)\\ (34.4,76.7)\\ (33.2,67.3) \end{array}$	$\begin{array}{c} (84.1,\ 79.8)\\ (\underline{89.8},\ 86.3)\\ (81.6,\ 74.4) \end{array}$	$\begin{array}{c}(71.4,\ 72.4)\\(78.3,\ 82.7)\\(69.1,\ 70.5)\end{array}$	(78.4, 76.4) (71.5, 78.9) (59.0, 67.2)	$\begin{array}{c} (\textbf{93.6},  \textbf{86.4}) \\ (89.7,  \underline{89.5}) \\ (\underline{95.6},  \underline{91.0}) \end{array}$	$\begin{array}{c} (\underline{91.4},\ \underline{85.2})\\ (94.8,\ 91.2)\\ (97.7,\ 92.4)\end{array}$	
	Average Rank	(38.1, 70.7) (6.0, 5.7)	$\begin{array}{c} (85.2,  80.2) \\ (2.7,  3.0) \end{array}$	(72.9, 75.2) (4.3, 4.3)	$\substack{(69.7,\ 74.1)\\(4.7,\ 5.0)}$	$\frac{(93.0, 89.0)}{(2.0, 1.7)}$	(94.6, 89.6) (1.3, 1.3)	
Pixel-level (AUROC, max-F1)	ISIC ColonDB ClinicDB TN3K	$\begin{array}{c} (83.8,\ 74.2)\\ (71.8,\ 31.5)\\ (66.2,\ 29.1)\\ (66.8,\ 32.6)\end{array}$	$\begin{array}{c} (67.1,\ 48.5)\\ (61.1,\ 19.6)\\ (67.1,\ 24.4)\\ (67.2,\ 30.0) \end{array}$	$\begin{array}{c} (\textbf{94.2},  \underline{79.6}) \\ (87.3,  \underline{56.5}) \\ (83.3,  \underline{56.2}) \\ (73.3,  35.7) \end{array}$	$\begin{array}{c} (\underline{94.2},  81.0) \\ (\overline{86.1},  45.7) \\ (\underline{83.5},  43.0) \\ (70.1,  32.5) \end{array}$	$\begin{array}{c} (92.1,\ 77.4)\\ (\underline{88.7},\ 52.6)\\ (82.5,\ 51.8)\\ (\underline{75.9},\ \underline{36.4})\end{array}$	$\begin{array}{c} (89.3,\ 71.4)\\ (\textbf{90.4},\ \textbf{58.2})\\ (\textbf{84.4},\ \textbf{58.2})\\ (\textbf{77.2},\ \textbf{41.9})\end{array}$	
	Average Rank	(72.1, 41.8) (5.5, 4.5)	(65.6, 30.6) (5.5, 6.0)	$\begin{array}{c} (84.5,  \underline{57.0}) \\ (\underline{2.5},  \underline{2.3}) \end{array}$	$\begin{array}{c} (83.5,\ 50.5)\\ (3.0,\ 3.5) \end{array}$	$\begin{array}{c} (\underline{84.8}, \ 54.6) \\ (2.8, \ 2.8) \end{array}$	(85.3, 57.4) (1.8, 2.0)	

the highest average rankings, showcasing robust generalization capabilities across different domains. As depicted in Fig. 3, AdaCLIP demonstrates precise detection of various anomalies across diverse medical categories, such as identifying skin cancer regions in photographic images and detecting thyroid nodules in ultrasound images. AdaCLIP achieves notably superior performance in locating abnormal lesion/tumor regions compared to other ZSAD methods. More quantitative and qualitative results in Appendix Section 5-7 further illustrate the superior ZSAD performance of AdaCLIP.



Fig. 3: Visualization of anomaly maps of different ZSAD methods. The proposed AdaCLIP can get the most precise segmentation results for novel categories in both industrial and medical domains.

## 5.3 Ablation Study

Influence of Prompts. Table 3 presents the detection performance of AdaCLIP with different combinations of static prompts and dynamic prompts, namely V1  $(w/o \mathbf{P}^S, w/o \mathbf{P}^D)$ , V2  $(w/\mathbf{P}^S, w/o \mathbf{P}^D)$ , V3  $(w/o \mathbf{P}^S, w/\mathbf{P}^D)$ , and V4  $(w/\mathbf{P}^S, w/\mathbf{P}^D)$ . The superior performance of V2 and V3 to V1 shows that both prompts are useful. V4 with hybrid prompts brings the most significant improvements. This is because static prompts struggle to capture diverse anomalies, while solely dynamic prompts are not sufficient. The combined hybrid prompts offer robust and flexible adaptation, thereby offering better ZSAD performance. Fig. 4 visualizes the patch embeddings and anomaly maps to delve into the influence of prompts. It clearly shows that both prompts are useful in highlighting the abnormal patch embeddings, facilitating precise predictions. However, with solely static or dynamic prompts, the prediction results are not perfect. In comparison, the model (V4) with hybrid prompts can detect anomalies more accurately. We also evaluate the influence of multimodal prompts and find it crucial to prompt both the text and image encoders, as shown in Appendix Section 2.1.

Analysis on Prompting Depth and Length. Fig. 5 visualizes the ZSAD performance of AdaCLIP under different prompting depths (J) and prompting lengths (M). Significantly, the performance of AdaCLIP does not exhibit continuous improvement with larger J and M. This is because the incorporation of more learnable prompting parameters introduces a risk of overfitting the auxiliary



Fig. 4: Visualization of Patch Embeddings and Anomaly Maps under Different Prompts. PCA is utilized to reduce the dimension of patch embeddings for enhanced visualization. For individual models, the left shows patch embeddings and the right displays anomaly maps.

**Table 3:** Ablation Results of Static prompts  $\mathbf{P}^{S}$  and Dynamic prompts  $\mathbf{P}^{D}$ .

Model	$del \mathbf{P}^{S} \mathbf{P}^{D}$		Medical	Domain	Industrial Domain		
			Image-level	Pixel-level	Image-level	Pixel-level	
V1	X	Х	(87.9, 58.3)	(83.9, 54.3)	(86.7, 85.0)	(92.8, 45.9)	
V2	~	X	(88.8, 60.4)	(84.8, 56.4)	(89.1, 88.7)	(94.1, 48.2)	
V3	X	~	(88.4, 60.0)	(84.4, 57.0)	(86.9, 87.1)	(93.5, 46.1)	
V4	~	~	(94.6, 89.6)	( <b>85.3</b> , <b>57.4</b> )	(90.2, 89.6)	(94.2, 50.2)	



 Table 4: Ablation results on HSF.

HSF	Medical	Domain	Industrial Domain		
	Image-level	Pixel-level	Image-level	Pixel-level	
×	(91.8, 88.7)	(85.7, 57.7)	(86.0, 85.8)	(93.8, 49.9	
~	(94.6, 89.6)	(85.3, 57.4)	(90.2, 89.6)	(94.2, 50.2)	



training dataset. To mitigate this, we employ the default setting J = 4 and M = 5, ensuring consistently high AD performance across both domains.

Influence of HSF. Table 4 showcases the impact of HSF. The results reveal a significant improvement of image-level ZSAD performance across both medical and industrial domains with the introduction of HSF compared to maximum-based image-level AD (without HSF). For instance, the image-level AUROC increases from 86.0% to 90.2% in the industrial domain. This improvement is attributed to the ability of HSF to aggregate the semantics of abnormal regions from multiple hierarchies. Conversely, relying on the maximum value of anomaly maps for image-level AD yields suboptimal results. Additional analysis of HSF is provided in Appendix Section 2.2.

**Influence of Annotated Auxiliary Data.** We conduct experiments in the medical domain to explore the influence of annotated auxiliary data, as illustrated in Table 5 and Fig. 6. Relying exclusively on medical datasets for training results

**Table 5:** ZSAD Performance in the Med-ical Domain with Varied Training Data.Top: Image-level AD. Bottom: Pixel-level AD.

Dataset	Medical	Industrial	Both
HeadCT	(76.0, 72.3)	(81.6, 78.8)	( <b>91.4</b> , <b>85.2</b> )
BrainMRI	(57.6, 76.0)	(86.4, 85.3)	(94.8, 91.2)
Br35H	(68.7, 68.9)	(68.8, 69.4)	(97.7, 92.4)
Average	(67.4, 72.4)	(78.9, 77.8)	(94.6, 89.6)
ISIC	(68.5, 49.5)	(89.2, 72.3)	(89.3, 71.4)
ColonDB	(89.4, 55.4)	(78.5, 31.3)	(90.4, 58.2)
ClinicDB	(91.3,  65.1)	(78.2, 39.5)	(84.4, 58.2)
TN3K	(69.7, 33.6)	(75.9, 41.4)	(77.2, <b>41.9</b> )
Average	(79.7, 50.9)	(80.5, 46.1)	( <b>85.3</b> , <b>57.4</b> )



Fig. 6: Anomaly Maps Visualization Across Different Training Sets. Categories and corresponding datasets for individual samples are listed on the left.

in subpar ZSAD performance, as illustrated by the notable underperformance on ISIC when trained solely with medical data. This can be attributed to the lack of data diversity within the selected medical dataset, such as ColonDB [45] or ClinicDB [4]. The utilized industrial datasets offer more diverse anomalies, thereby providing greater generalization capacity when trained with them. Notably, training with ColonDB brings surprisingly promising results on ClinicDB, even surpassing more diverse training sets. This is because ColonDB and ClinicDB both focus on colon polyp detection and thus, these two datasets share similarities despite being acquired through different imaging techniques, as shown in Fig. 6. Generally, using more diverse auxiliary training sets can improve the generalization ability.

Influence of Backbones. Table 6 illustrates the impact of different backbones. AdaCLIP demonstrates significantly improved results in both medical and industrial domains with a larger backbone, ViT-L/14@336px, compared to ViT-B/16. Moreover, the additional parameters are lightweight compared to the original CLIP parameters, comprising only 4.6% (40.7 MB) of the original parameters added to ViT-L/14@336px (890.8 MB). This effectively demonstrates that existing VLMs can be adapted to ZSAD using lightweight parameters.

#### 5.4 Analysis

**Rationale behind the ZSAD Scheme with Auxiliary Data.** The ZSAD scheme with auxiliary data successfully tailors existing VLMs, including DINOV2, SAM, and CLIP, for ZSAD. To explore the reason why training with auxiliary data can improve ZSAD performance, we visually analyze the distributions of patch embeddings from these models across two datasets featuring diverse anomalies, *i.e.*, MVTec and VisA. In Fig. 7, it becomes evident that abnormal patch embeddings in both MVTec and VisA exhibit distinctive characteristics compared to the normal ones. Meanwhile, the normal embeddings in these two datasets exhibit similar distributions. Consequently, the decision boundary learned in MVTec is applicable to VisA despite not being trained on VisA. This

**Table 6:** Comparison between various back-bones. Sizes of original CLIP and added pa-rameters by AdaCLIP parameters are re-ported in Mega Bytes.

В	ackbone	ViT-B/16	ViT-L/14@336px
Size	(Ori., Added)	(334.6, 19.6)	(890.8, 40.7)
Industria	l Image-level	$\begin{array}{c} (81.3,\ 78.3) \\ (82.5,\ 52.7) \end{array}$	(94.6, 89.6)
Domain	Pixel-level		(85.3, 57.4)
Medical	Image-level	$\begin{array}{c} (83.9,\ 84.6)\\ (91.7,\ 42.1) \end{array}$	(90.2, 89.6)
Domain	Pixel-level		(94.2, 50.2)



**Fig. 7:** t-SNE [36] visualization of normal/abnormal patch embeddings.

phenomenon can be attributed to the high-level similarities in normalities and abnormalities present in both datasets as perceived by VLMs. The awareness of these similarities can be harnessed by learning annotated auxiliary data.

Enhancing ZSAD Performance through Prompt Optimization. The influence of prompting tokens on predictions becomes apparent in both Table 3 and Fig. 4. While prompts generated by AdaCLIP are promising, the potential for further improving ZSAD performance exists through prompt optimization. We leverage model V4 in Table 3 and refine its prompts for specific images using corresponding anomaly masks for training. The results are depicted in the right two columns of Fig. 4, illustrating that optimized prompts result in more discernible abnormal patch embeddings and finer anomaly maps, particularly noticeable in the bottom two rows. This underscores the significance of devising methods to generate optimal prompts tailored to individual images.

# 6 Conclusion

In this study, we introduce AdaCLIP, a generic ZSAD model to detect anomalies across arbitrary novel categories without any reference image. AdaCLIP leverages annotated auxiliary AD data for training and effectively adapts pre-trained CLIP for ZSAD by integrating learnable hybrid prompts. Additionally, a HSF module is proposed to extract region-level anomaly information to enhance imagelevel AD performance. Through extensive experimentation across 14 datasets spanning industrial and medical domains, AdaCLIP demonstrates promising AD performance in novel categories from different domains.

**Discussion and Limitations.** Our experimental results demonstrate the potential of AdaCLIP as a powerful solution for ZSAD. We believe that leveraging more diverse annotated auxiliary anomaly detection datasets can improve the generalization capability of AdaCLIP. In fact, like any other ZSAD method, AdaCLIP might fail when testing data that significantly depart from auxiliary training data, as shown in Sec. 7.1 in the Appendix.

# Acknowledgements

This paper is supported in part by FAIR (Future Artificial Intelligence Research) project, funded by the NextGenerationEU program within the PNRR-PE-AI scheme (M4C2, Investment 1.3, Line on Artificial Intelligence), in part by the Ministry of Industry and Information Technology of the People's Republic of China under Grant #2023ZY01089, and in part by the China Scholarship Council (CSC) under Grant 202306160078.

#### References

- Aota, T., Tong, L.T.T., Okatani, T.: Zero-shot versus many-shot: Unsupervised texture anomaly detection. In: 2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV). pp. 5553–5561 (2023)
- Arthur, D., Vassilvitskii, S., et al.: k-means++: The advantages of careful seeding. In: Soda. vol. 7, pp. 1027–1035 (2007)
- Bergmann, P., Batzner, K., Fauser, M., Sattlegger, D., Steger, C.: The MVTec anomaly detection dataset: A comprehensive real-world dataset for unsupervised anomaly detection. International Journal of Computer Vision 129(4), 1038–1059 (2021)
- Bernal, J., Sánchez, F.J., Fernández-Esparrach, G., Gil, D., Rodríguez, C., Vilariño, F.: Wm-dova maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from physicians. Computerized medical imaging and graphics 43, 99–111 (2015)
- Bionda, A., Frittoli, L., Boracchi, G.: Deep autoencoders for anomaly detection in textured images using CW-SSIM. In: International Conference on Image Analysis and Processing (ICIAP). pp. 669–680 (2022)
- Cai, Y., Liang, D., Luo, D., He, X., Yang, X., Bai, X.: A discrepancy aware framework for robust anomaly detection. IEEE Transactions on Industrial Informatics pp. 1–10 (2023)
- Çallı, E., Sogancioglu, E., van Ginneken, B., van Leeuwen, K.G., Murphy, K.: Deep learning for chest x-ray analysis: A survey. Medical Image Analysis 72, 102125 (2021)
- Cao, Y., Wan, Q., Shen, W., Gao, L.: Informative knowledge distillation for image anomaly segmentation. Knowledge-Based Systems 248, 108846 (2022)
- Cao, Y., Xu, X., Liu, Z., Shen, W.: Collaborative discrepancy optimization for reliable image anomaly localization. IEEE Transactions on Industrial Informatics pp. 1–10 (2023)
- Cao, Y., Xu, X., Sun, C., Cheng, Y., Du, Z., Gao, L., Shen, W.: Segment any anomaly without training via hybrid prompt regularization. arXiv preprint arXiv:2305.10724 (2023)
- Cao, Y., Xu, X., Sun, C., Gao, L., Shen, W.: Bias: Incorporating biased knowledge to boost unsupervised image anomaly localization. IEEE Transactions on Systems, Man, and Cybernetics: Systems (2023)
- Cao, Y., Xu, X., Zhang, J., Cheng, Y., Huang, X., Pang, G., Shen, W.: A survey on visual anomaly detection: Challenge, approach, and prospect. arXiv preprint arXiv:2401.16402 (2024)
- Chandola, V., Banerjee, A., Kumar, V.: Anomaly detection: A survey. ACM Computing Surveys 41(3) (2009)

- 16 Y. Cao et al.
- Chen, X., Han, Y., Zhang, J.: 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)
- Chen, X., Zhang, J., Tian, G., He, H., Zhang, W., Wang, Y., Wang, C., Wu, Y., Liu, Y.: Clip-ad: A language-guided staged dual-path model for zero-shot anomaly detection. arXiv preprint arXiv:2311.00453 (2023)
- Deng, H., Li, X.: Anomaly detection via reverse distillation from one-class embedding. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9737–9746 (2022)
- Ding, C., Pang, G., Shen, C.: Catching both gray and black swans: Open-set supervised anomaly detection. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7388–7398 (2022)
- Gong, H., Chen, G., Wang, R., Xie, X., Mao, M., Yu, Y., Chen, F., Li, G.: Multi-task learning for thyroid nodule segmentation with thyroid region prior. In: 2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI). pp. 257–261 (2021)
- 19. Gu, Z., Zhu, B., Zhu, G., Chen, Y., Tang, M., Wang, J.: Anomalygpt: Detecting industrial anomalies using large vision-language models. In: Proceedings of the AAAI conference on artificial intelligence (2024)
- 20. Gutman, D., Codella, N.C.F., Celebi, E., Helba, B., Marchetti, M., Mishra, N., Halpern, A.: Skin lesion analysis toward melanoma detection: A challenge at the international symposium on biomedical imaging (isbi) 2016, hosted by the international skin imaging collaboration (isic). arXiv preprint arXiv:1605.01397 (2016)
- Hamada, A.: Br35h: Brain tumor detection 2020 (2020), https://www.kaggle.com/ datasets/ahmedhamada0/brain-tumor-detection
- He, H., Zhang, J., Chen, H., Chen, X., Li, Z., Chen, X., Wang, Y., Wang, C., Xie, L.: Diad: A diffusion-based framework for multi-class anomaly detection. In: AAAI (2024)
- Huang, C., Guan, H., Jiang, A., Zhang, Y., Spratlin, M., Wang, Y.: Registration based few-shot anomaly detection. In: European Conference on Computer Vision (2022)
- Jeong, J., Zou, Y., Kim, T., Zhang, D., Ravichandran, A., Dabeer, O.: Winclip: Zero-/few-shot anomaly classification and segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 19606–19616 (June 2023)
- Jezek, S., Jonak, M., Burget, R., Dvorak, P., Skotak, M.: Deep learning-based defect detection of metal parts: evaluating current methods in complex conditions. In: 2021 13th International Congress on Ultra Modern Telecommunications and Control Systems and Workshops (ICUMT). pp. 66–71 (2021)
- Jia, M., Tang, L., Chen, B.C., Cardie, C., Belongie, S., Hariharan, B., Lim, S.N.: Visual prompt tuning. In: European Conference on Computer Vision. pp. 709–727. Springer (2022)
- Jiang, Y., Cao, Y., Shen, W.: A masked reverse knowledge distillation method incorporating global and local information for image anomaly detection. Knowledge-Based Systems p. 110982 (2023)
- Kanade, P.B., Gumaste, P.: Brain tumor detection using mri images. Brain 3(2), 146–150 (2015)
- Khattak, M.U., Rasheed, H., Maaz, M., Khan, S., Khan, F.S.: Maple: Multi-modal prompt learning. In: The IEEE/CVF Conference on Computer Vision and Pattern Recognition (2023)

- Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., Dollar, P., Girshick, R.: Segment anything. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 4015–4026 (October 2023)
- 31. Kitamura, F.C.: Head ct hemorrhage (2018). https://doi.org/10.34740/KAGGLE/ DSV/152137, https://www.kaggle.com/dsv/152137
- Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: Proceedings of the IEEE international conference on computer vision. pp. 2980–2988 (2017)
- 33. Liu, J., Xie, G., Chen, R., Li, X., Wang, J., Liu, Y., Wang, C., Zheng, F.: Real3d-ad: A dataset of point cloud anomaly detection. arXiv preprint arXiv:2309.13226 (2023)
- Liu, M., Jiao, Y., Lu, J., Chen, H.: Anomaly detection for medical images using teacher-student model with skip connections and multiscale anomaly consistency. IEEE Transactions on Instrumentation and Measurement 73, 1–15 (2024). https: //doi.org/10.1109/TIM.2024.3406792
- 35. Liu, S., Zeng, Z., Ren, T., Li, F., Zhang, H., Yang, J., Li, C., Yang, J., Su, H., Zhu, J., et al.: Grounding dino: Marrying dino with grounded pre-training for open-set object detection. arXiv preprint arXiv:2303.05499 (2023)
- Van der Maaten, L., Hinton, G.: Visualizing data using t-sne. Journal of machine learning research 9(11) (2008)
- Milletari, F., Navab, N., Ahmadi, S.A.: V-net: Fully convolutional neural networks for volumetric medical image segmentation. In: 2016 fourth international conference on 3D vision (3DV). pp. 565–571. Ieee (2016)
- Mishra, P., Verk, R., Fornasier, D., Piciarelli, C., Foresti, G.L.: Vt-adl: A vision transformer network for image anomaly detection and localization. In: 2021 IEEE 30th International Symposium on Industrial Electronics (ISIE). pp. 01–06. IEEE (2021)
- Oquab, M., Darcet, T., Moutakanni, T., Vo, H.V., Szafraniec, M., Khalidov, V., Fernandez, P., Haziza, D., Massa, F., El-Nouby, A., Howes, R., Huang, P.Y., Xu, H., Sharma, V., Li, S.W., Galuba, W., Rabbat, M., Assran, M., Ballas, N., Synnaeve, G., Misra, I., Jegou, H., Mairal, J., Labatut, P., Joulin, A., Bojanowski, P.: Dinov2: Learning robust visual features without supervision. arXiv:2304.07193 (2023)
- Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021)
- Roth, K., Pemula, L., Zepeda, J., Schölkopf, B., Brox, T., Gehler, P.: Towards total recall in industrial anomaly detection. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 14318–14328 (2022)
- Ruff, L., Vandermeulen, R.A., Görnitz, N., Binder, A., Müller, E., Müller, K.R., Kloft, M.: Deep semi-supervised anomaly detection. In: International Conference on Learning Representations (2020)
- Schuhmann, C., Kaczmarczyk, R., Komatsuzaki, A., Katta, A., Vencu, R., Beaumont, R., Jitsev, J., Coombes, T., Mullis, C.: Laion-400m: Open dataset of clipfiltered 400 million image-text pairs. In: NeurIPS Workshop Datacentric AI. Jülich Supercomputing Center (2021)
- Tabernik, D., Sela, S., Skvarc, J., Skoaj, D.: Segmentation-based deep-learning approach for surface-defect detection. Journal of Intelligent Manufacturing **31**, 759 – 776 (2019)

- 18 Y. Cao et al.
- Tajbakhsh, N., Gurudu, S.R., Liang, J.: Automated polyp detection in colonoscopy videos using shape and context information. IEEE Transactions on Medical Imaging 35(2), 630–644 (2016)
- Tamura, M.: Random word data augmentation with clip for zero-shot anomaly detection. In: 34th British Machine Vision Conference 2023, BMVC 2023, Aberdeen, UK, November 20-24, 2023. BMVA (2023), https://papers.bmvc2023.org/0018. pdf
- Wan, Q., Gao, L., Li, X., Wen, L.: Industrial image anomaly localization based on gaussian clustering of pretrained feature. IEEE Transactions on Industrial Electronics 69(6), 6182–6192 (2022)
- 48. Wang, C., Zhu, W., Gao, B.B., Gan, Z., Zhang, J., Gu, Z., Qian, S., Chen, M., Ma, L.: Real-iad: A real-world multi-view dataset for benchmarking versatile industrial anomaly detection. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 22883–22892 (2024)
- Wieler, M., Hahn, T.: Weakly supervised learning for industrial optical inspection. In: DAGM symposium in. vol. 6 (2007)
- Yao, H., Yu, W., Wang, X.: A Feature Memory Rearrangement Network for Visual Inspection of Textured Surface Defects Toward Edge Intelligent Manufacturing. IEEE Transactions on Automation Science and Engineering pp. 1–20 (2022)
- 51. Zavrtanik, V., Kristan, M., Skočaj, D.: Dsr a dual subspace re-projection network for surface anomaly detection. In: European Conference on Computer Vision (2022)
- Zha, Y., Wang, J., Dai, T., Chen, B., Wang, Z., Xia, S.T.: Instance-aware dynamic prompt tuning for pre-trained point cloud models. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (2023)
- Zhang, J., Chen, X., Wang, Y., Wang, C., Liu, Y., Li, X., Yang, M.H., Tao, D.: Exploring plain vit reconstruction for multi-class unsupervised anomaly detection. arXiv preprint arXiv:2312.07495 (2023)
- Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Conditional prompt learning for visionlanguage models. In: 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 16795–16804 (2022)
- 55. Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Learning to prompt for vision-language models. International Journal of Computer Vision **130**(9), 2337–2348 (2022)
- Zhou, Q., Pang, G., Tian, Y., He, S., Chen, J.: Anomalyclip: Object-agnostic prompt learning for zero-shot anomaly detection. In: International Conference on Learning Representations (2024)
- 57. Zou, Y., Jeong, J., Pemula, L., Zhang, D., Dabeer, O.: Spot-the-difference selfsupervised pre-training for anomaly detection and segmentation. In: European Conference on Computer Vision. pp. 392–408. Springer (2022)