

# AdaptCLIP: Adapting CLIP for Universal Visual Anomaly Detection

Bin-Bin Gao<sup>1</sup> Yue Zhou<sup>2,3</sup> Jiangtao Yan<sup>1</sup> Yuezhi Cai<sup>2</sup>  
 Weixi Zhang<sup>2</sup> Meng Wang<sup>2</sup> Jun Liu<sup>1</sup> Yong Liu<sup>1</sup> Lei Wang<sup>2</sup> Chengjie Wang<sup>1,4</sup>  
<sup>1</sup>Tencent YouTu Lab <sup>2</sup>Siemens AG <sup>3</sup>Technical University of Munich <sup>4</sup>Shanghai Jiao Tong University

## Abstract

Universal visual anomaly detection aims to identify anomalies from novel or unseen vision domains without additional fine-tuning, which is critical in open scenarios. Recent studies have demonstrated that pre-trained vision-language models like CLIP exhibit strong generalization with just zero or a few normal images. However, existing methods struggle with designing prompt templates, complex token interactions, or requiring additional fine-tuning, resulting in limited flexibility. In this work, we present a simple yet effective method called AdaptCLIP based on two key insights. First, adaptive visual and textual representations should be learned alternately rather than jointly. Second, comparative learning between query and normal image prompt should incorporate both contextual and aligned residual features, rather than relying solely on residual features. AdaptCLIP treats CLIP models as a foundational service, adding only three simple adapters, visual adapter, textual adapter, and prompt-query adapter, at its input or output ends. AdaptCLIP supports zero-/few-shot generalization across domains and possesses a training-free manner on target domains once trained on a base dataset. AdaptCLIP achieves state-of-the-art performance on 12 anomaly detection benchmarks from industrial and medical domains, significantly outperforming existing competitive methods. We will make the code and model of AdaptCLIP available at <https://github.com/gaobb/AdaptCLIP>.

## 1. Introduction

Universal visual anomaly detection (AD) aims to identify anomaly images and segment anomaly pixels from novel or unseen visual objects after learning a single model on a base or seen dataset. This is a more challenging task as it requires strong generalization when facing cross-domain datasets. Meanwhile, it is a more practical topic as people are more interested in fast adaptability in real-world scenarios, especially in low data regimes (i.e., few-shot and even zero-shot). For example, in medical image diagnosis and industrial visual quality inspection, it is difficult to

collect a large-scale dataset due to inherent scarcity and privacy protection. Recently, developing universal visual AD has attracted increasing attention because existing unsupervised ADs with either separated [8, 25, 34] or unified models [11, 48] perform poorly in unseen objects despite promising performance on seen objects.

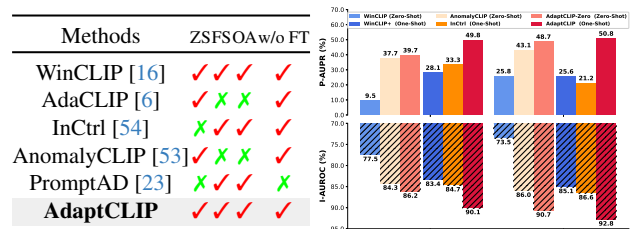


Figure 1. Comparisons of state-of-the-arts and our AdaptCLIP. ✓ means satisfied and ✗ means not satisfied. Our method supports zero-/few-shot (ZS and FS) visual AD across different domains without fine-tuning (FT) on the target dataset. It only adds simple adapters at CLIP’s input or output ends without complex token interactions, thus preserving CLIP’s original ability (OA). The AdaptCLIP using only one normal image prompt achieves the best performance in image-level anomaly classification (I-AUROC) and pixel-level anomaly segmentation (P-AUPR) on 12 AD benchmarks from industrial and medical domains. Moreover, the zero-shot AdaptCLIP is also significantly better than existing zero-shot and even some one-shot approaches. The detailed results are reported in Tabs. 1 and 2. Best viewed in zoom.

To address this fragmentation, recent works have attempted to design universal models to recognize anomalies for unseen objects. They typically build on vision-language models (i.e., CLIP [29]) benefiting from strong generalization. WinCLIP [16] computes anomaly scores on dense patch windows. This brings large computational costs and memory burden, limiting high-resolution input or large pre-trained models. AnomalyCLIP [53] learns class-agnostic prompt embeddings to align patch-wise tokens thus avoiding dense window operations. It further refines vanilla CLIP by concatenating learnable tokens to intermediate layers of CLIP. AdaCLIP [6] further integrates visual knowledge into textual prompt embeddings. However, they may destroy inherent representations of CLIP. Therefore, *we want to explore whether we can achieve the same or even better AD*

performance while maintaining the original ability.

In contrast, humans perceive anomalies when an input significantly deviates from those normal patterns stored in our brains. There is evidence to support this point in neuroscience [31]. PatchCore [34] builds a memory bank storing normal features and PaDiM [7] learns a multivariate Gaussian distribution of normal features. At inference, anomalies are recognized by comparing input features with the memory bank or the learned distribution. However, these methods usually require a certain number of normal images and thus are limited in universal (i.e., open-world) scenarios. Two recent works, i.e., InCtrl [54] and PromptAD [23], have studied how to further improve performance with few-shot normal images. However, InCtrl only considers anomaly classification, while PromptAD needs to learn a new model for each class. Different from them, *we want to comprehensively explore a universal AD model, aiming to detect any anomalies in image-level and pixel-level from cross-domains without any training on target domains.*

Toward this end, we propose a simple but effective universal visual anomaly detection framework, called AdaptCLIP. *The philosophy of AdaptCLIP is that “less and simpler could be better”, and it contains three adapters designed by two key insights :* First, adaptive visual and textual representations should be learned alternately rather than jointly. Second, comparative learning between query and the corresponding normal image prompt should incorporate both contextual and aligned residual features, rather than relying solely on residual features. Our contributions are summarized as follows.

- We propose a simple but effective universal visual anomaly detection framework based on visual-language CLIP models, which is capable of detecting any visual anomalies at image- and pixel-level from cross-domain datasets without any training on target domains.
- We propose visual and textual adapters, and find that they should alternately learn adaptive visual and textual representation guided by the powerful vision-language representations from CLIP models.
- We propose a prompt-query adapter that aims to capture meta-perceptual capabilities between query image and the corresponding normal image prompt, based on their joint distribution of contextual features of the query and the aligned residual features between prompt and query.
- AdaptCLIP outperforms zero- and few-shot AD methods on 8 industrial and 4 medical benchmarks, as shown in Fig. 5. Meanwhile, AdaptCLIP possesses simpler adapters, fewer parameters, and competitive efficiency.

## 2. Related Works

**Unsupervised ADs** target to identify anomalies given sufficient normal training images. Most unsupervised

AD methods can be roughly grouped into three categories: embedding-, discrimination-, and reconstruction-based methods. Embedding-based methods, such as PaDiM [7], MDND [32], PatchCore [34], CS-Flow [35] and PyramidFlow [19], assume that offline features extracted from a pre-trained model preserve discriminative information and thus help to separate anomalies from normal samples. Discrimination-based methods, such as Cut-Paste [21], DRAEM [50], and SimpleNet [25], typically convert unsupervised AD to supervised ones by introducing pseudo (synthesized) anomaly samples. Reconstruction-based ADs, such as autoencoder [2, 12, 14, 40], generative adversarial networks [27, 45, 49] and reconstruction networks [24, 33, 51], assume that anomalous regions should not be able to be properly reconstructed and thus result in high reconstruction errors since they do not exist in normal training samples. The recent knowledge distillation [4, 8, 36, 42, 43] or feature reconstruction methods [11, 46, 48, 52] train a student or reconstruction network to match a fixed pre-trained teacher network and achieve a good balance between effectiveness and efficiency. However, all these methods are limited to recognizing anomalies of seen classes but often perform poorly on unseen classes. For a novel scenario, people have to collect sufficient normal images first and then retrain a model. This is inefficient and lacks the rapid adaptability required for practical applications.

**Zero-Shot ADs** have achieved impressive performance by utilizing large vision-language models, e.g., CLIP [29]. WinCLIP [16] designs two-class textual prompts and introduces multi-scale patch windows for accurate anomaly segmentation. It brings large computational costs and memory burden, limiting high-resolution input or large pre-trained models. AnomalyCLIP [53] learns class-agnostic prompt embeddings to align patch-wise tokens thus avoiding dense window operation. In addition, AnomalyCLIP refines vanilla CLIP representation by appending some learnable tokens to the middle layer of CLIP. Recently, AdaCLIP [6] and VCP-CLIP [28] utilize similar ideas and further integrate visual knowledge into textual prompt embeddings. We argue that these additional operations make models more complex and may hurt the original capabilities of CLIP. Instead of visual-language models, ACR [20] and MuSc [22] perform zero-shot AD only requiring batch-level and full-shot testing images, but they may be limited in privacy protection scenarios. Different from these methods, we explore whether the same or even better AD performance is achieved while retaining the original ability of CLIP without any information on test data distribution.

**Few-Shot ADs** mainly pay attention to learning or using only a limited number of normal images, such as TDG [37], RegAD [15], GraphCore [44] and FastRecon [10]. Some works [9, 47] consider another few-shot setting where a

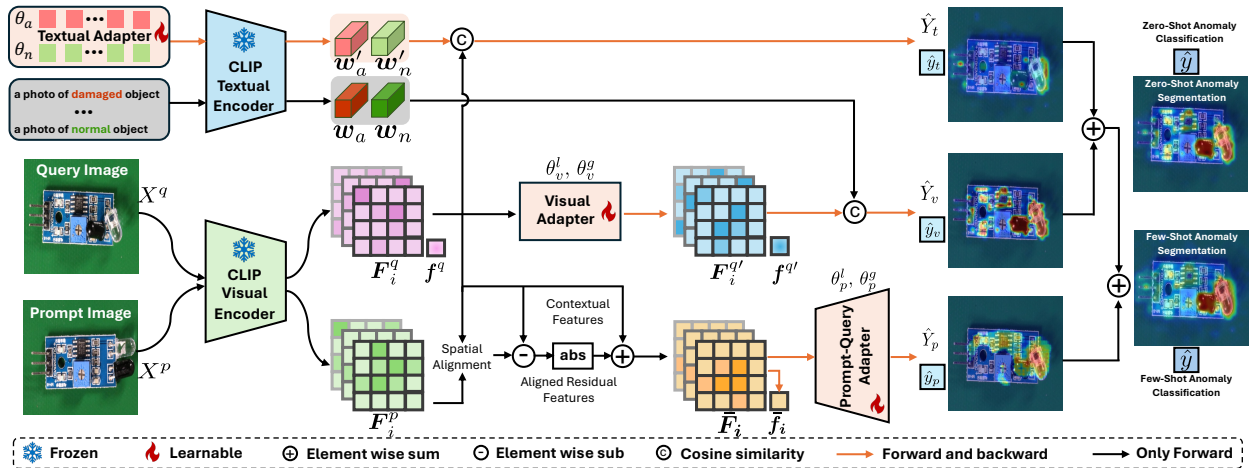


Figure 2. The framework of **AdaptCLIP**, which consists of three pluggable adapters, i.e., visual adapter, textual adapter, and prompt-query adapter. First, the first two adapters alternately learn visual and textual representations for zero-shot anomaly detection (Sec. 3.2). The prompt-query adapter further learns a comparison ability between query image and its corresponding normal prompt for few-shot anomaly detection (Sec. 3.3). Once trained, it can segment any anomalies providing only few-shot and even zero-shot normal image prompts.

limited number of samples is given from anomaly images. The performance of these methods lags behind unsupervised ADs. Recently, few-shot AD performance has been improved significantly by visual-language models. WinCLIP+ [16] is the first work to apply CLIP models to few-shot AD, which stores normal tokens into a memory bank, then retrieves the nearest token for each query token using cosine similarity, and finally computes an anomaly map using the nearest distance. InCtrl [54] further integrates multi-level information, including patch-level residual maps and image-level residual features, and prior knowledge score using two-class textual prompts, to learn a holistic scoring function for anomaly classification. It does not consider pixel-level anomaly segmentation. PromptAD [23] introduces the concept of explicit anomaly margin, which mitigates the training challenge caused by the absence of anomaly training images. However, it requires re-training models when applied to target datasets. In contrast, we explore jointly optimizing anomaly classification and segmentation in a unified model, which can quickly adapt to novel scenarios only given few-shot normal image prompts, not involving additional re-training.

### 3. Methods

**Problem Formulation:** Our objective is to learn a universal AD model that detects any anomalies from diverse domains without any training on target dataset. Thus, a reasonable assumption is that there is a different distribution between training and testing sets. Formally, let  $\mathcal{D}_{\text{base}} = \{X_i, Y_i, y_i\}_{i=1}^N$  be a training dataset, that consists of  $N$  normal and anomalous images,  $X_i \in \mathcal{R}^{h \times w \times 3}$  is the  $i$ -th image, and  $Y_i \in \mathcal{R}^{h \times w}$  and  $y_i = \{0, 1\}$  is the corresponding anomaly mask and anomaly label, with  $y_i = 0$  indicates normal and  $y_i = 1$  signifies anomaly. The testing set  $\mathcal{T}$

may consist of multiple different domains with various objects and anomaly types. Here, we denote the  $t$ -th novel domain as  $\mathcal{D}_{\text{novel}}^t = \{X_i, Y_i, y_i\}_{i=1}^{N_t}$ . Under a few-shot setting, a few normal images  $\mathcal{P}_c = \{X_i\}_{i=1}^k$  are randomly drawn from each class of the target domain, where  $c$  is the class index and  $k$  is typically a small number, e.g.,  $k = \{1, 2, 4\}$ . It is worth noting that  $\mathcal{P}_c$  is only available during inference, and cannot be used in any way during training phase.

**Overview:** As illustrated in Fig. 2, the visual adapter adapts patch and image tokens with fixed two-class textual prompt embeddings. The textual adapter learns two-class prompt embeddings to align with the fixed patch and image tokens. The prompt-query adapter operates in a one-prompt meta-learning manner, leveraging the joint distribution of query context features and the aligned residual features between the prompt and query. In a zero-shot scenario, image-level anomaly score and pixel-level anomaly map can be obtained using textual and visual adapters (Sec. 3.2). In a few-shot scenario, anomaly score and map are derived by integrating predictions from zero-shot and prompt-query adapters (Sec. 3.3). Below we present them in detail.

#### 3.1. Revisiting CLIP for Anomaly Detection

For a query image  $X^q \in \mathcal{R}^{h \times w \times 3}$ , we feed it to visual encoder  $\mathcal{F}(\cdot)$  and obtain local patch tokens  $\{F_i^q \in \mathcal{R}^d\}_{i=1}^{hw/p^2}$  and global image token  $f^q \in \mathcal{R}^d$ , where  $p$  is patch size. WinCLIP [16] introduces two-class prompts describing normal and abnormal states. For example, “a photo of a normal object” and “a photo of a damaged object”. In practical application, one could design multiple textual descriptions for normal and abnormal states. Feeding these normal and abnormal descriptions to textual encoder  $\mathcal{T}(\cdot)$ , we can obtain the embeddings of normal  $w_n \in \mathcal{R}^d$  and abnormal  $w_a \in \mathcal{R}^d$ . The pixel-level anomaly map is computed by

measuring the cosine similarities between all patch tokens and the textual embeddings, that is

$$\hat{Y} = \left[ \frac{\exp(\langle \mathbf{w}_a, \mathbf{F}_i^q \rangle)}{\exp(\langle \mathbf{w}_a, \mathbf{F}_i^q \rangle) + \exp(\langle \mathbf{w}_n, \mathbf{F}_i^q \rangle)} \right], \quad (1)$$

where  $\langle \cdot \rangle$  represents the cosine similarity, and  $[\cdot]$  means that all patch-wise prediction scores are rearranged according to their spatial positions and interpolated to the original input resolution. Replacing  $\mathbf{F}_i^q$  with  $\mathbf{f}^q$  in Eq. 1, we can obtain an image-level anomaly score  $\hat{y}$  for  $X^q$ , that is

$$\hat{y} = \frac{\exp(\langle \mathbf{w}_a, \mathbf{f}^q \rangle)}{\exp(\langle \mathbf{w}_a, \mathbf{f}^q \rangle) + \exp(\langle \mathbf{w}_n, \mathbf{f}^q \rangle)}. \quad (2)$$

### 3.2. AdaptCLIP with Alternating Learning

To adapt CLIP for universal visual anomaly detection, we design visual and textual adapters to alternately learn visual and textual representations. Specifically, the visual adapter learns adaptive visual tokens ( $\mathbf{F}_i^{q'}$  and  $\mathbf{f}^{q'}$ ) when fixing two-class static textual embeddings ( $\mathbf{w}_a$  and  $\mathbf{w}_n$ ), while the textual adapter learns two-class textual prompt embeddings ( $\mathbf{w}'_a$  and  $\mathbf{w}'_n$ ) when fixing visual tokens ( $\mathbf{F}_i^q$  and  $\mathbf{f}^q$ ).

**Visual Adapter** adapts vision tokens ( $\mathbf{F}_i^q$  and  $\mathbf{f}^q$ ) with fixed textual embeddings ( $\mathbf{w}_a$  and  $\mathbf{w}_n$ ). It consists of two branches, global and local, which transform global image token and local patch tokens, respectively. Architecturally, the global and local branches are implemented using a simple residual multi-layer perception (MLP), that is

$$\mathbf{F}_i^{q'} = \mathbf{F}_i^q + \text{MLP}(\mathbf{F}_i^q; \theta_v^l); \mathbf{f}^{q'} = \mathbf{f}^q + \text{MLP}(\mathbf{f}^q; \theta_v^g), \quad (3)$$

where  $\theta_v^l$  and  $\theta_v^g$  are learnable parameters. Replacing  $\mathbf{F}_i^q$  and  $\mathbf{f}^q$  in Eqs. 1 and 2 with  $\mathbf{F}_i^{q'}$  and  $\mathbf{f}^{q'}$ , we obtain pixel-level anomaly map  $\hat{Y}_v$  and image-level anomaly score  $\hat{y}_v$ .

**Textual Adapter** aims to directly learn two-class prompts  $\theta_a, \theta_n \in \mathcal{R}^{r \times d}$  without prompt templates, where  $r > 0$  is the length of prompts. We feed them into the frozen textual encoder  $\mathcal{T}(\cdot)$  of CLIP, and obtain the corresponding embeddings  $\mathbf{w}'_a$  and  $\mathbf{w}'_n$ , that is

$$\mathbf{w}'_a = \mathcal{T}(\theta_a), \mathbf{w}'_n = \mathcal{T}(\theta_n). \quad (4)$$

Then, we replace the static  $\mathbf{w}_a$  and  $\mathbf{w}_n$  in Eqs. 1 and 2 with the learnable prompt embeddings  $\mathbf{w}'_a$  and  $\mathbf{w}'_n$  to derive local and global anomaly predictions,  $\hat{Y}_t$  and  $\hat{y}_t$ .

**Alternating Learning or Joint Learning?** A possible question is whether we can learn visual and textual representations jointly. That is, in Eqs. 1 and 2, we simultaneously replace fixed textual embeddings and visual tokens with learnable prompt embeddings ( $\mathbf{w}'_a$  and  $\mathbf{w}'_n$ ) and adaptive visual tokens ( $\mathbf{F}_i^{q'}$  and  $\mathbf{f}^{q'}$ ). Indeed, this joint alignment mechanism is successful when a large-scale image-text dataset is available. However, we empirically find that it does not work well in the AD field, as shown in Tab. 5

(Lines 3 vs. 4). This is not surprising because the available training data scale is still relatively small and lacks fine-grained textual annotations. The joint learning easily overfits and leads to poor generalization on novel datasets. In contrast, the alternating learning helps us fully utilize the prior knowledge of the CLIP model and thus improve the cross-domain generalization.

### 3.3. AdaptCLIP with Comparative Learning

Compared to static or learnable textual prompts, using a normal image as a visual prompt is more intuitive. Therefore, we expect to learn a comparison ability between a query image  $X^q$  and its corresponding normal prompt  $X^p$ , which generalizes well to unseen objects. We find that applying multi-layer features yields better results. For simplicity, we use a single-layer feature in the following.

**Spatial Alignment:** A simple way is to directly measure their difference by the absolute value of their residual feature, that is  $|\mathbf{F}_i^q - \mathbf{F}_i^p|$ , where  $\mathbf{F}_i^q$  and  $\mathbf{F}_i^p$  are the patch token of  $X^q$  and  $X^p$ , respectively. It may fail if the query and prompt images are not aligned in pixel space (e.g., due to rotation and translation). Therefore, we have to align query and prompt tokens for effective comparison. For any query token  $\mathbf{F}_i^q$ , we search the nearest one among all normal tokens  $\{\mathbf{F}_j^p\}_{j=1}^{hw/p^2}$  using euclidean distance, that is

$$\mathbf{F}_i^{p'} = \mathbf{F}_k^p, k = \arg \min_j \|\mathbf{F}_i^q - \mathbf{F}_j^p\|_2. \quad (5)$$

Then, we take  $\mathbf{F}_i^{p'}$  as aligned prompt token of  $\mathbf{F}_i^q$ . Now, we can derive the aligned residual feature, i.e.,  $|\mathbf{F}_i^q - \mathbf{F}_i^{p'}|$ .

**Joint contextual and aligned residual feature:** The aligned residual feature highlights differences or anomaly regions well. However, it may lose contextual information or introduce noise. Intuitively, the contextual information is critical to identify anomalies. Therefore, we aggregate the original query tokens and the aligned residual features by an element-wise sum,

$$\bar{\mathbf{F}}_i = \mathbf{F}_i^q + |\mathbf{F}_i^q - \mathbf{F}_i^{p'}|. \quad (6)$$

**Prompt-Query Adapter:** The ultimate goal is to achieve pixel-level anomaly segmentation and image-level anomaly classification. Therefore, we propose a lightweight segmentation head  $\mathcal{G}(\cdot; \theta_p^l)$  to learn anomaly segmentation based on the joint feature  $\bar{\mathbf{F}}$ , that is

$$\hat{Y}_p = \mathcal{G}(\bar{\mathbf{F}}; \theta_p^l), \quad (7)$$

where  $\theta_p^l$  is its parameters. Specifically, the segmentation head consists of several transposed convolution blocks following a  $1 \times 1$  convolution layer. Here, each transposed convolution block upsamples input feature by  $2 \times$ , and it is composed of a  $3 \times 3$  convolution, a BatchNorm, a ReLU, and a  $2 \times 2$  deconvolution.

Meanwhile, we need to obtain a global image-level prediction. First, we perform average-pooling and max-pooling on the joint feature  $\bar{F}$  along the spatial dimension and then take their weighted average as the global image representation. Then, a simple MLP is used to map the global feature to an image-level prediction score, that is

$$\hat{y}_p = \text{MLP}((\text{AvgPool}(\bar{F}) + \text{MaxPool}(\bar{F}))/2; \theta_p^g), \quad (8)$$

where  $\theta_p^g$  is the parameter.

---

```

"""
image_encoder and textual_encoder: frozen CLIP model
visual_adapter: learned visual adapter
theta: learned two-class prompt embeddings
prompt_query_adapter: learned prompt-query adapter
Xq, Xp: query image and the normal image prompt
T: fixed two-class textual prompts
inference_style: 'zero-shot' or 'few-shot' (default)
"""
# extract fixed visual and textual representations
Fq = image_encoder(Xq), W = textual_encoder(T)
# adapt visual tokens and textual embeddings, Eqs. 3 and 4
Fq_a = visual_adapter(Fq), W_a = textual_encoder(theta)

# AdaptCLIP via alternating learning for zero-shot AD
# compute anomaly map with softmax, Eq. 1
Yv = F.softmax(F.cosine_similarity(W, Fq_a))
Yt = F.softmax(F.cosine_similarity(W_a, Fq))
Y = (Yv + Yt) / 2

# AdaptCLIP via comparative learning for few-shot AD
if inference_style == 'few-shot':
    Fp = image_encoder(Xp)
    # spatial alignment, Eq. 5
    Fp_a = Fp[torch.min((Fq - Fp)**2)[1], dim = -1)]
    # joint contextual and aligned residual feature, Eq. 6
    Fq_ = Fq + torch.abs(Fq - Fp_a)
    # perform anomaly segmentation, Eq. 7
    Yp = prompt_query_adapter(Fq_)
    Y = (Yv + Yt + Yp) / 3

```

---

Figure 3. PyTorch pseudocode for the inference of AdaptCLIP.

### 3.4. Training and Inference

During training, we use cross-entropy loss for global image anomaly classification, and Focal and Dice losses for local patch anomaly segmentation, which is exactly the same as AnomalyCLIP [53]. For zero-shot inference, we average the predictions from visual and textual adapters. For few-shot inference, we fuse (i.e., average) all results from three adapters, i.e., prompt-query, visual and textual adapters, as the final predictions of AdaptCLIP. Fig. 3 shows PyTorch pseudocode for the pixel-level inference of AdaptCLIP, where  $W = [w_a^T; w_n^T]$ ,  $W_a = [w_a'^T; w_n'^T]$ ,  $Fq$  and  $Fp$  means the fixed  $F_i^q$  and  $F_i^p$ , and  $Fq_a$  and  $Fp_a$  refers the adaptive  $F^{q'}$  and the aligned  $F^{p'}$ . Here, we omit the image-level inference since it can be easily obtained by replacing local patch tokens with a global image token.

## 4. Experiments

### 4.1. Experimental Setup

**Datasets:** We comprehensively evaluate AdaptCLIP on multiple datasets from industrial and medical domains.

For industrial domain, we use MVTEC [3], VisA [55], BTAD [26], MVTEC3D [5], DTD [1], KSDD [38], MPDD [17], and large-scale Real-IAD [41]. In medical domain, we utilize brain tumor detection datasets, Br35H [13] and COVID-19 [30], as well as gastrointestinal polyp datasets, Kvasir [18] and Endo [39]. A detailed introduction to these datasets can be found in Appendix.

**Evaluation Metrics:** Following previous works, we use AUROC for image-level anomaly classification and AUPR for pixel-level anomaly segmentation in our main paper. Here, we emphasize that AUPR is better for anomaly segmentation, where the imbalance issue is very extreme between normal and anomaly pixels [55]. In Appendix, we also provide detailed comparisons using all metrics, including AUROC, AUPR, and  $F1_{\max}$ .

**Training and Testing Protocol:** Following AnomalyCLIP [53], we train AdaptCLIP using the testing data from MVTEC and evaluate zero-/few-shot performance on other datasets. As for the evaluation of MVTEC, we train AdaptCLIP using the testing data of VisA. For fair comparison, all models are trained and evaluated using the same protocol.

**Competing Methods:** We compare our AdaptCLIP with diverse state-of-the-art zero-/few-shot AD methods including zero-shot WinCLIP [16], AnomalyCLIP [53], AdaCLIP [6], and few-shot WinCLIP+ [16], InCtrl [54] and AnomalyCLIP+. Here, AnomalyCLIP+ is a strong baseline we build on AnomalyCLIP [53] by adding patch-level feature associations like WinCLIP+. More implementation details about AdaptCLIP and competing methods can be found in Appendix.

### 4.2. Comparisons with Zero-/Few-Shot Methods

Tab. 1 and 2 present comparisons of AdaptCLIP to competing zero-/few-shot methods in image-level anomaly classification and pixel-level anomaly segmentation, respectively, on 8 real-world industrial and 4 medical AD datasets. Note that we only use image-level metrics to evaluate Br35H and Covid due to the lack of pixel-level annotations, and only report the results for Kvasir and Endo using pixel-level metrics since normal images are not included in these two datasets. Below we analyze these results in detail.

**Generalization on Industrial Domain:** Generally, AdaptCLIP significantly outperforms all competing models on almost all industrial datasets across three few-shot settings, 1-shot, 2-shot and 4-shot. The performance of all methods generally gets better with more image prompts. Specifically, InCtrl [54] surpasses WinCLIP [16] due to additional fine-tuning on a base training dataset. AnomalyCLIP [53] further achieves better generalization, which verifies the importance of learning object-agnostic prompts. AdaptCLIP exhibits superior performance, outperforming AnomalyCLIP [53] by a large margin (about 10%+ in pixel AUPR and 2%+ in image AUROC), particularly on chal-

Table 1. **Image-level anomaly classification comparisons with AUROC metric on industrial and medical domains.** The best and second-best results are highlighted in red and blue, respectively. The superscript<sup>†</sup> indicates that the results are our re-implementation with the same training and testing protocol as AnomalyCLIP and our AdaptCLIP. Note that the results are averaged over all categories on each dataset and the full results of each category are presented in Appendix, the same below.

Shots	Methods	Industrial									Medical		
		MVTec	VisA	BTAD	MVTec3D	DTD	KSSD	MPDD	Real-IAD	AVG	Br35H	Covid	AVG
0	WinCLIP [16]	90.4	75.5	68.2	69.4	95.1	92.9	61.5	67.0	77.5	80.5	66.4	73.5
	AdaCLIP <sup>†</sup> [6]	90.7	81.7	89.9	76.2	92.7	96.6	64.0	73.3	83.1	96.7	69.4	83.0
	AnomalyCLIP [53]	91.6	82.0	88.3	73.9	93.9	97.8	77.5	69.5	84.3	94.2	77.7	86.0
	AdaptCLIP-Zero	93.5	84.8	91.0	78.6	96.0	98.1	73.6	74.2	86.2	94.8	86.5	90.7
1	WinCLIP+ [16]	93.6±0.4	80.0±2.4	84.4±1.5	74.1±0.4	97.9±0.2	93.8±0.4	69.3±2.9	74.7±0.2	83.4	80.1±2.1	90.1±3.6	85.1
	InCtrl [54]	91.3±0.4	83.2±2.4	88.5±0.4	75.3±1.3	97.9±0.3	92.0±0.9	73.0±2.7	76.6±0.0	84.7	83.9±6.4	89.2±5.3	86.6
	AnomalyCLIP+ [53]	95.2±0.2	86.1±0.7	88.5±0.8	76.7±2.1	98.0±0.2	97.5±0.3	83.4±2.6	78.2±0.0	88.0	90.8±5.1	87.3±2.6	89.1
	AdaptCLIP	94.5±0.5	90.5±1.2	93.4±0.0	81.7±1.5	98.0±0.0	96.9±0.3	83.8±2.2	81.8±0.3	90.1	93.7±2.4	91.8±2.5	92.8
2	WinCLIP+ [16]	94.5±1.0	82.7±1.0	85.8±1.8	74.3±0.3	98.1±0.2	93.8±0.2	69.3±2.3	76.1±0.1	84.3	81.6±0.6	91.8±2.5	86.7
	InCtrl [54]	91.8±0.9	86.3±1.4	86.2±2.0	75.4±0.5	98.3±0.2	91.6±0.9	74.2±1.8	78.5±0.0	85.3	86.1±1.7	89.7±5.1	87.9
	AnomalyCLIP+ [53]	95.4±0.1	87.8±0.5	89.2±1.1	78.3±1.3	98.2±0.1	97.9±0.2	83.4±1.5	78.3±0.0	88.6	91.5±4.0	89.3±2.7	90.4
	AdaptCLIP	95.7±0.6	92.2±0.8	93.4±0.2	82.9±1.1	98.3±0.0	97.2±0.0	84.4±0.7	82.9±0.2	90.8	94.0±1.7	94.9±0.9	94.5
4	WinCLIP+ [16]	95.3±0.1	84.3±0.6	87.8±0.8	75.7±0.3	98.2±0.0	94.0±0.2	71.2±1.6	77.0±0.0	85.4	82.3±0.4	92.9±2.1	87.6
	InCtrl [54]	93.1±0.7	87.8±0.2	67.5±2.4	78.1±1.1	97.7±0.1	91.6±0.9	78.6±2.3	81.8±0.0	84.5	89.1±1.2	91.4±4.1	90.3
	AnomalyCLIP+ [53]	96.1±0.1	88.8±0.5	90.5±1.2	79.2±1.3	98.4±0.1	97.8±0.1	86.3±1.8	78.4±0.0	89.4	91.1±4.4	91.4±3.0	91.3
	AdaptCLIP	96.6±0.3	93.1±0.2	93.3±0.3	84.2±0.6	98.5±0.1	97.0±0.2	86.8±1.1	83.9±0.2	91.7	93.7±2.0	95.8±0.9	94.8

Table 2. **Pixel-level anomaly segmentation comparisons with AUPR metric on industrial and medical domains.**

Shots	Methods	Industrial									Medical		
		MVTec	VisA	BTAD	MVTec3D	DTD	KSSD	MPDD	Real-IAD	AVG	Kvasir	Endo	AVG
0	WinCLIP [16]	18.2	5.4	12.9	5.3	9.8	7.1	14.1	3.3	9.5	27.8	23.8	25.8
	AdaCLIP <sup>†</sup> [6]	39.1	31.0	42.9	37.5	75.2	48.2	25.9	30.5	41.3	36.6	43.7	40.1
	AnomalyCLIP [53]	34.5	21.3	45.5	30.5	62.6	51.9	28.9	26.7	37.7	39.6	46.6	43.1
	AdaptCLIP-Zero	38.3	26.1	41.8	31.4	68.7	58.3	25.3	28.2	39.7	45.3	52.0	48.7
1	WinCLIP+ [16]	38.3±0.8	15.8±0.2	41.3±2.6	18.4±1.1	47.8±0.9	19.2±0.3	29.8±2.0	13.9±0.2	28.1	27.6±2.9	23.6±0.1	25.6
	InCtrl [54]	47.8±1.1	17.7±0.6	44.1±1.4	18.7±0.5	64.3±0.5	26.7±0.7	27.9±2.2	19.1±0.0	33.3	22.1±1.7	20.3±3.7	21.2
	AnomalyCLIP+ [53]	40.8±0.1	24.8±0.9	41.3±1.1	30.6±1.1	67.4±0.4	47.5±0.5	34.2±0.8	27.9±0.0	39.3	46.9±3.9	47.8±4.9	47.4
	AdaptCLIP	53.7±0.9	38.9±0.3	60.6±1.0	40.7±0.6	76.9±0.1	57.8±1.2	33.5±2.5	36.6±0.1	49.8	49.2±4.7	52.4±4.7	50.8
2	WinCLIP+ [16]	39.5±0.6	17.2±0.8	42.8±1.3	19.1±0.8	48.2±0.9	19.0±0.5	30.7±1.1	14.8±0.1	28.9	29.1±0.2	27.6±2.3	28.4
	InCtrl [54]	49.2±0.7	18.5±0.2	44.2±0.8	20.3±0.6	64.4±0.4	26.4±2.5	29.2±1.3	20.1±0.0	34.0	24.9±1.9	24.5±7.5	24.7
	AnomalyCLIP+ [53]	41.5±0.1	26.2±0.7	41.9±0.6	32.4±1.5	68.1±0.2	47.6±0.4	35.3±1.1	28.1±0.0	40.1	47.3±2.9	49.6±4.8	48.5
	AdaptCLIP	55.1±0.5	40.7±0.6	61.0±0.6	42.3±1.1	77.4±0.2	57.5±1.1	35.0±0.7	37.8±0.1	50.9	49.0±4.1	53.1±4.2	51.1
4	WinCLIP+ [16]	41.2±0.9	18.1±1.3	44.0±0.4	19.9±0.6	49.3±0.1	19.1±0.7	32.0±0.2	15.4±0.2	29.9	29.6±0.8	27.7±0.5	28.7
	InCtrl [54]	50.9±0.3	19.2±0.6	44.0±0.2	22.2±1.2	64.9±0.3	26.0±1.4	31.4±0.8	21.0±0.0	35.0	24.7±1.6	22.3±1.0	23.5
	AnomalyCLIP+ [53]	42.4±0.0	27.5±1.1	45.8±3.0	33.4±1.3	68.5±0.2	46.4±0.7	36.8±1.0	28.2±0.0	41.1	45.9±1.5	49.2±3.4	47.6
	AdaptCLIP	57.2±0.8	41.8±0.6	62.3±0.3	44.5±0.3	78.2±0.2	56.4±1.4	37.4±1.1	39.1±0.3	52.1	47.5±2.7	52.2±3.1	49.9

lenging and large-scale datasets like VisA and Real-IAD. This reveals the power of comparative learning based on the joint contextual and aligned residual features for universal anomaly detection. Under zero-shot setting, AdaptCLIP-Zero significantly outperforms SOTA AdaCLIP [6] on anomaly classification, although it shows a slight weakness in industrial anomaly segmentation. However, AdaptCLIP is simpler, requires fewer learnable parameters (0.6M vs. 10.7M in Tab. 3), and generalizes better from the industrial to the medical domain. In addition, our one-shot AdaptCLIP easily outperforms zero-shot AdaCLIP [6] if only one

normal image prompt is available.

**Generalization on Medical Domain:** Our AdaptCLIP performs strongly on medical AD regardless of zero-shot or few-shot settings when applying the same model trained on an industrial dataset (i.e., MVTEC). Surprisingly, it significantly outperforms SOTA AdaCLIP on image anomaly classification (i.e., 6.3% in AUROC) and pixel anomaly segmentation (i.e., 8.6% in AUPR). Notably, our approach still works even when replacing normal image prompts with anomaly images. This is meaningful for some special datasets that don’t contain any normal images, such as

Table 3. Complexity and efficiency comparisons.

Shots	Methods	CLIP Models	Input Size	# Params (M)	Inf.Time (ms)
0	WinCLIP [16]	ViT-B-16+240	240×240	208.4 + 0.0	201.3
		ViT-B-16+240	512×512	208.4 + 0.0	3912.6
	AdaCLIP [6]	ViT-L/14@336px	518×518	428.8 + 10.7	212.0
		ViT-L/14@336px	518×518	427.9 + 5.6	154.9
	AnomalyCLIP [53]	ViT-B-16+240	512×512	208.4 + 0.4	49.9
		ViT-L/14@336px	518×518	427.9 + 0.6	162.2
1	WinCLIP+ [16]	ViT-B-16+240	240×240	208.4 + 0.0	339.5
		ViT-B-16+240	512×512	208.4 + 0.0	7434.9
	InCtrl [54]	ViT-B-16+240	240×240	208.4 + 0.3	337.0
		ViT-B-16+240	518×518	427.9 + 5.6	158.6
	AnomalyCLIP+ [53]	ViT-B-16+240	512×512	208.4 + 1.4	54.0
		ViT-L/14@336px	518×518	427.9 + 1.8	168.2

Kvasir and Endo. Here, this success is mainly due to the proposed spatial alignment mechanism, as well as a strong prior assumption that anomaly pixels are mostly sparse.

**Efficiency Comparison:** We measure complexity and efficiency by the number of parameters and the forward inference time, as shown in Tab. 3. The evaluation is performed on one V100 GPU with batch size 32. The number of parameters of AdaCLIP and AnomalyCLIP is 17 times and 9 times that of our AdaptCLIP, respectively. Compared to SOTA, AdaptCLIP achieves competitive inference time yet better AD performance. When extending from zero-shot to one-shot, AnomalyCLIP+ and our AdaptCLIP require almost no additional inference time, unlike earlier WinCLIP.

**Qualitative Results:** Fig. 4 shows some selected visualizations from industrial and medical testing images using AdaptCLIP. Generally, few-shot normal image prompts help AdaptCLIP segment anomalies more accurately and produce fewer false positives than in a zero-shot manner.

### 4.3. Comparisons with Many-Shot Methods

In Tab. 4, we compare few-shot AdaptCLIP with many-shot and full-shot unified AD models. It can be seen that AdaptCLIP is better than the early many-shot methods, RegAD [15], and comparable to the latest PromptAD [23]. It is worth noting that PromptAD [23] requires re-training with few-shot normal images while our method remains training-free on target domains. Furthermore, our method outperforms full-shot methods, such as SimpleNet [25] and UniAD [48], and is also competitive with the latest OneNIP [11]. In short, our method has shown excellent performance, especially in the open-world scenario for universal anomaly detection, although there is still some small gap compared to state-of-the-art full-shot methods.

### 4.4. Ablation Studies

To demonstrate the effectiveness of the proposed three adapters in AdaptCLIP, TA: Textual Adapter, VA: Visual Adapter, and PQA: Prompt-Query Adapter, and two main insights, alternating learning, and comparative learning based on the joint contextual and aligned residual feature, we conduct experiments on MVTec and VisA, and report results in Tab. 5.

Table 4. Comparisons of image-level anomaly classification and pixel-level anomaly segmentation (using AUROC/AUPR metric, and the same as below) with many-shot and full-shot methods.

Methods	Shots	MVTec	VisA	BTAD
<b>AdaptCLIP</b>	1	94.5 / 53.7	90.5 / 38.9	<b>93.4 / 60.6</b>
	4	<b>96.6 / 57.2</b>	<b>93.1 / 41.8</b>	<b>93.3 / 62.3</b>
RegAD [15]	8	91.2 / 51.1	79.7 / 28.6	90.7 / 40.5
PromptAD [23]	4	<b>96.6 / 52.9</b>	89.1 / 31.5	-
SimpleNet [25]	full	78.2 / 24.8	89.2 / 33.1	90.3 / 36.2
UniAD [48]	full	96.5 / 44.7	90.8 / 33.6	92.2 / 50.9
OneNIP [11]	full	<b>97.9 / 63.7</b>	<b>92.5 / 43.3</b>	92.6 / 56.8

Table 5. Ablation studies about different components.

No.	Methods	Shots	TA	VA	PQA	MVTec	VisA
0		0	✗	✗	✗	91.1 / 33.0	82.1 / 18.0
1	baselines	0	✓	✗	✗	92.2 / 31.4	82.9 / 19.7
2		0	✗	✓	✗	90.5 / 39.4	81.0 / 22.1
3	joint	0	✓	✓	✗	89.3 / 36.2	81.6 / 21.5
4	<b>alternating</b>	0	✓	✓	✗	93.5 / 38.3	84.8 / 26.1
5	w/o context	1	✗	✗	✓	62.6 / 7.0	85.3 / 28.7
6	<b>w context</b>	1	✗	✗	✓	88.1 / 50.2	88.9 / 38.1
7	<b>AdaptCLIP</b>	1	✓	✓	✓	<b>94.2 / 52.5</b>	<b>92.0 / 38.8</b>

**Simple but effective baselines.** The first baseline is the naive CLIP (Line 0), and it is simple and effective for zero-shot anomaly detection only using two-class textual prompts. However, it is still weak in pixel-level anomaly segmentation. The individual textual adapter and visual adapter are two additional baselines. Specifically, the textual adapter can be seen as an extreme simplification of AnomalyCLIP [53], removing the textual prompt template and textual prompt tuning. The simple textual adapter performs better than the original AnomalyCLIP and naive CLIP in anomaly classification, although it is slightly inferior in anomaly segmentation (Lines 0 vs. 1). The visual adapter learns adaptive local patch tokens and global image tokens to align textual representations from CLIP in both patch and image levels. This significantly improves pixel-level anomaly segmentation (Lines 0 vs. 2).

**Alternating learning is better than joint learning.** We explore the impact of alternating learning and joint learning strategies on AdaptCLIP’s performance. Alternating learning adapts visual or textual representations independently, while joint learning optimizes both representations simultaneously. As shown (Lines 3 vs. 4) in Tab. 5, the alternating learning strategy significantly enhances the performance of AdaptCLIP compared to joint learning. Alternating learning not only fully leverages the strong prior guidance of CLIP’s visual and textual representations but also mitigates the risk of over-fitting due to fine-tuning on a small training dataset. Additionally, we observe that the visual adapter

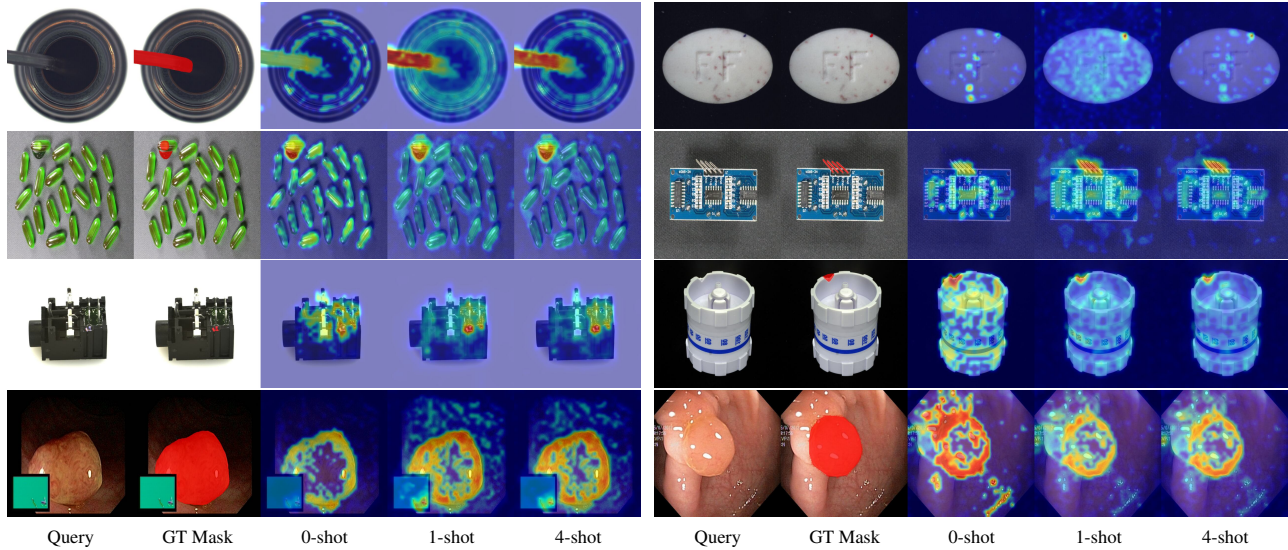


Figure 4. Qualitative comparisons of our AdaptCLIP with different prompt numbers on MVTec, VisA, Real-IAD, Kvasir and Endo. More qualitative results of AdaptCLIP can be found in Appendix. Best viewed in color and zoom.

alone excels in anomaly segmentation (Line 2), whereas the textual adapter alone performs better in anomaly classification (Line 1). By integrating the alternating learning into visual and textual adapters, AdaptCLIP generally achieves superior anomaly detection performance (Line 4).

**The joint of contextual information and aligned residual features performs better than residual features alone.**

The aligned residual feature captures the distinctions between anomalous features and their corresponding normal counterparts. It effectively eliminates features related to individual objects and may improve generalization. However, we realize that isolated residual features may lose contextual information about visual objects, resulting in degraded model performance or even training failure (Line 5). Therefore, we propose a joint feature learning based on both contextual and aligned residual features, which further significantly boosts the model’s performance (Lines 6 vs. 5). This means contextual information is equally important for anomaly identification. Notably, the optimal performance for AdaptCLIP is achieved when all proposed components are integrated (Line 7).

**Effects on pre-trained CLIP models.** We report zero- and one-shot results of AdaptCLIP using different CLIP models in Tab. 6. It can be seen that a larger pre-trained model always brings better performance, especially in image-level classification. Furthermore, our method equipped with a lightweight model (ViT-B-16+240) makes it possible to achieve competitive anomaly segmentation performance.

**5. Conclusion**

In this paper, we introduce a universal anomaly detection task, which focuses on generalizing anomaly detection models across domains, such as industrial and medical, and

Table 6. Ablation studies about different pre-trained CLIP models.

CLIP Models	Input Size	Shots	MVTec	VisA
ViT-B-16+240	512×512	0	83.9 / 38.3	75.4 / 19.5
ViT-L/14@336px	518×518	0	93.5 / 38.3	84.8 / 26.1
ViT-B-16+240	512×512	1	92.4 / 52.3	85.2 / 30.3
ViT-L/14@336px	518×518	1	94.2 / 52.5	92.0 / 38.8

in open scenarios, such as zero- or few-shot settings. Once the universal anomaly detection model is trained, it does not need any fine-tuning on the target dataset. Compared with single zero-shot or few-shot AD models, the universal anomaly detection model is more flexible, supporting zero-/few-shot inference via fixed or learnable textual prompts and a few normal image prompts, while providing both image-level and pixel-level anomaly predictions. We propose a universal anomaly detection framework, AdaptCLIP, which alternately learns adaptive visual representations and text prompt embeddings, as well as jointly learns comparisons based on the contextual information of query image and the aligned residual features between the query and the prompt. Extensive experiments on 8 standard industrial and 4 medical datasets show that AdaptCLIP significantly outperforms current competitive models in multiple settings.

**Limitation:** AdaptCLIP achieves good AD performance only given zero-/few-shot normal image prompts. However, it could cause the model to confuse normal and abnormal instances and finally result in a decreased performance when we provide anomaly images as normal image prompts. Fortunately, normal images are generally relatively easy to obtain in practical applications. In addition, it may work using abnormal images as visual prompts because most of the pixels may be normal even in anomaly images.

## References

- [1] Toshimichi Aota, Lloyd Teh Tzer Tong, and Takayuki Okatani. Zero-shot versus many-shot: Unsupervised texture anomaly detection. In *WACV*, 2023. 5
- [2] Yoshua Bengio, Li Yao, Guillaume Alain, and Pascal Vincent. Generalized denoising auto-encoders as generative models. In *NeurIPS*, 2013. 2
- [3] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. MVTEC-AD: A comprehensive real-world dataset for unsupervised anomaly detection. In *CVPR*, 2019. 5
- [4] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. Uninformed Students: Student-teacher anomaly detection with discriminative latent embeddings. In *CVPR*, 2020. 2
- [5] Paul Bergmann., Xin Jin., David Sattlegger., and Carsten Steger. The mvtec 3d-ad dataset for unsupervised 3d anomaly detection and localization. In *VISAPP*, 2022. 5
- [6] 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 *ECCV*, 2024. 1, 2, 5, 6, 7, 4
- [7] Thomas Defard, Aleksandr Setkov, Angelique Loesch, and Romaric Audigier. PaDiM: a patch distribution modeling framework for anomaly detection and localization. In *ICPR*, 2021. 2
- [8] Hanqiu Deng and Xingyu Li. Anomaly detection via reverse distillation from one-class embedding. In *CVPR*, 2022. 1, 2
- [9] Choubo Ding, Guansong Pang, and Chunhua Shen. Catching both gray and black swans: Open-set supervised anomaly detection. In *CVPR*, 2022. 2
- [10] Zheng Fang, Xiaoyang Wang, Haocheng Li, Jiejie Liu, Qiguai Hu, and Jimin Xiao. FastRecon: Few-shot industrial anomaly detection via fast feature reconstruction. In *ICCV*, 2023. 2
- [11] Bin-Bin Gao. Learning to detect multi-class anomalies with just one normal image prompt. In *ECCV*, 2024. 1, 2, 7
- [12] Dong Gong, Lingqiao Liu, Vuong Le, Budhaditya Saha, Moussa Reda Mansour, Svetha Venkatesh, and Anton van den Hengel. Memorizing normality to detect anomaly: Memory-augmented deep autoencoder for unsupervised anomaly detection. In *ICCV*, 2019. 2
- [13] Ahmed Hamada. Br35h: Brain tumor detection 2020, 2020. 5
- [14] Jinlei Hou, Yingying Zhang, Qiaoyong Zhong, Di Xie, Shiliang Pu, and Hong Zhou. Divide-and-Assemble: Learning block-wise memory for unsupervised anomaly detection. In *ICCV*, 2021. 2
- [15] Chaoqin Huang, Haoyan Guan, Aofan Jiang, Ya Zhang, Michael Spratling, and Yan-Feng Wang. Registration based few-shot anomaly detection. In *ECCV*, 2022. 2, 7
- [16] Jongheon Jeong, Yang Zou, Taewan Kim, Dongqing Zhang, Avinash Ravichandran, and Onkar Dabeer. WinCLIP: Zero-/few-shot anomaly classification and segmentation. In *CVPR*, 2023. 1, 2, 3, 5, 6, 7, 4
- [17] Stepan Jezek, Martin Jonak, Radim Burget, Pavel Dvorak, and Milos Skotak. Deep learning-based defect detection of metal parts: evaluating current methods in complex conditions. In *ICUMT*, 2021. 5
- [18] Debesh Jha, Pia H Smedsrud, Michael A Riegler, Pål Halvorsen, Thomas de Lange, Dag Johansen, and Håvard D Johansen. Kvasir-seg: A segmented polyp dataset. In *International Conference on Multimedia Modeling*, pages 451–462, 2020. 5
- [19] Jiarui Lei, Xiaobo Hu, Yue Wang, and Dong Liu. Pyramid-Flow: High-resolution defect contrastive localization using pyramid normalizing flow. In *CVPR*, 2023. 2
- [20] Aodong Li, Chen Qiu, Marius Kloft, Padhraic Smyth, Maja Rudolph, and Stephan Mandt. Zero-shot anomaly detection via batch normalization. In *NeurIPS*, 2023. 2
- [21] Chun-Liang Li, Kihyuk Sohn, Jinsung Yoon, and Tomas Pfister. CutPaste: Self-supervised learning for anomaly detection and localization. In *CVPR*, 2021. 2
- [22] Xurui Li, Ziming Huang, Feng Xue, and Yu Zhou. MuSc: Zero-shot industrial anomaly classification and segmentation with mutual scoring of the unlabeled images. In *ICLR*, 2024. 2
- [23] Xiaofan Li, Zhizhong Zhang, Xin Tan, Chengwei Chen, Yanyun Qu, Yuan Xie, and Lizhuang Ma. PromptAD: Learning prompts with only normal samples for few-shot anomaly detection. In *CVPR*, 2024. 1, 2, 3, 7
- [24] Wenrui Liu, Hong Chang, Bingpeng Ma, Shiguang Shan, and Xilin Chen. Diversity-measurable anomaly detection. In *CVPR*, 2023. 2
- [25] Zhikang Liu, Yiming Zhou, Yuansheng Xu, and Zilei Wang. SimpleNet: A simple network for image anomaly detection and localization. In *CVPR*, 2023. 1, 2, 7
- [26] Pankaj Mishra, Riccardo Verk, Daniele Fornasier, Claudio Picciarelli, and Gian Luca Foresti. VT-ADL: A vision transformer network for image anomaly detection and localization. In *ISIE*, 2021. 5
- [27] Pramuditha Perera, Ramesh Nallapati, and Bing Xiang. OCGAN: One-class novelty detection using GANs with constrained latent representations. In *CVPR*, 2019. 2
- [28] 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. In *ECCV*, 2024. 2
- [29] 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 *ICML*, 2021. 1, 2
- [30] Tawsifur Rahman, Amith Khandakar, Yazan Qiblawey, Anas Tahir, Serkan Kiranyaz, Saad Bin Abul Kashem, Mohammad Tariqul Islam, Somaya Al Maadeed, Susu M Zughaier, Muhammad Salman Khan, et al. Exploring the effect of image enhancement techniques on covid-19 detection using chest x-ray images. *Computers in biology and medicine*, 2021. 5
- [31] Rajesh PN Rao and Dana H Ballard. Predictive coding in the visual cortex: a functional interpretation of some extraclassical receptive-field effects. *Nature Neuroscience*, 2(1), 1999. 2

- [32] Oliver Rippel, Patrick Mertens, and Dorit Merhof. Modeling the distribution of normal data in pretrained deep features for anomaly detection. In *ICPR*, 2021. 2
- [33] Nicolae-Cătălin Ristea, Neelu Madan, Radu Tudor Ionescu, Kamal Nasrollahi, Fahad Shahbaz Khan, Thomas B Moeslund, and Mubarak Shah. Self-supervised predictive convolutional attentive block for anomaly detection. In *CVPR*, 2022. 2
- [34] Karsten Roth, Latha Pemula, Joaquin Zepeda, Bernhard Schölkopf, Thomas Brox, and Peter Gehler. Towards total recall in industrial anomaly detection. In *CVPR*, 2022. 1, 2
- [35] Marco Rudolph, Tom Wehrbein, Bodo Rosenhahn, and Bastian Wandt. Fully convolutional cross-scale-flows for image-based defect detection. In *WACV*, 2022. 2
- [36] Mohammadreza Salehi, Niousha Sadjadi, Soroosh Baselizadeh, Mohammad H Rohban, and Hamid R Rabiee. Multiresolution knowledge distillation for anomaly detection. In *CVPR*, 2021. 2
- [37] Shelly Sheynin, Sagie Benaim, and Lior Wolf. A hierarchical transformation-discriminating generative model for few shot anomaly detection. In *ICCV*, 2021. 2
- [38] Domen Tabernik, Samo Šela, Jure Skvarč, and Danijel Skočaj. Segmentation-based deep-learning approach for surface-defect detection. *Journal of Intelligent Manufacturing*, 31(3):759–776, 2020. 5
- [39] David Vázquez, Jorge Bernal, F Javier Sánchez, Gloria Fernández-Esparrach, Antonio M López, Adriana Romero, Michal Drozdal, Aaron Courville, et al. A benchmark for endoluminal scene segmentation of colonoscopy images. *Journal of healthcare engineering*, 2017, 2017. 5
- [40] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *ICML*, 2008. 2
- [41] Chengjie Wang, Wenbing Zhu, Bin-Bin Gao, Zhenye Gan, Jiangning Zhang, Zhihao Gu, Shuguang Qian, Mingang Chen, and Lizhuang Ma. Real-iad: A real-world multi-view dataset for benchmarking versatile industrial anomaly detection. In *CVPR*, 2024. 5
- [42] Guodong Wang, Shumin Han, Errui Ding, and Di Huang. Student-teacher feature pyramid matching for anomaly detection. *BMVC*, 2021. 2
- [43] Shenzi Wang, Liwei Wu, Lei Cui, and Yujun Shen. Glancing at the patch: Anomaly localization with global and local feature comparison. In *CVPR*, 2021. 2
- [44] Guoyang Xie, Jingbao Wang, Jiaqi Liu, Feng Zheng, and Yaochu Jin. Pushing the limits of fewshot anomaly detection in industry vision: Graphcore. In *ICLR*, 2023. 2
- [45] Xudong Yan, Huaidong Zhang, Xuemiao Xu, Xiaowei Hu, and Pheng-Ann Heng. Learning semantic context from normal samples for unsupervised anomaly detection. In *AAAI*, 2021. 2
- [46] Xincheng Yao, Ruoqi Li, Zefeng Qian, Yan Luo, and Chongyang Zhang. Focus the Discrepancy: Intra-and inter-correlation learning for image anomaly detection. In *ICCV*, 2023. 2
- [47] Xincheng Yao, Ruoqi Li, Jing Zhang, Jun Sun, and Chongyang Zhang. Explicit boundary guided semi-push-pull contrastive learning for supervised anomaly detection. In *CVPR*, 2023. 2
- [48] Zhiyuan You, Lei Cui, Yujun Shen, Kai Yang, Xin Lu, Yu Zheng, and Xinyi Le. A unified model for multi-class anomaly detection. In *NeurIPS*, 2022. 1, 2, 7
- [49] Muhammad Zaigham Zaheer, Jin-ha Lee, Marcella Astrid, and Seung-Ik Lee. Old is Gold: Redefining the adversarially learned one-class classifier training paradigm. In *CVPR*, 2020. 2
- [50] Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. DRAEM: A discriminatively trained reconstruction embedding for surface anomaly detection. In *ICCV*, 2021. 2
- [51] Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. Reconstruction by inpainting for visual anomaly detection. *PR*, 112, 2021. 2
- [52] Ying Zhao. OmniAL: A unified cnn framework for unsupervised anomaly localization. In *CVPR*, 2023. 2
- [53] Qihang Zhou, Guansong Pang, Yu Tian, Shibo He, and Jiming Chen. AnomalyCLIP: Object-agnostic prompt learning for zero-shot anomaly detection. In *ICLR*, 2024. 1, 2, 5, 6, 7, 4
- [54] Jiawen Zhu and Guansong Pang. Toward generalist anomaly detection via in-context residual learning with few-shot sample prompts. In *CVPR*, 2024. 1, 2, 3, 5, 6, 7, 4
- [55] Yang Zou, Jongheon Jeong, Latha Pemula, Dongqing Zhang, and Onkar Dabeer. Spot-the-difference self-supervised pre-training for anomaly detection and segmentation. In *ECCV*, 2022. 5