

# LaRP-CLIP: Layer-Aware Refinement with Prototype Guidance for Zero-Shot Anomaly Detection

Xing Fang<sup>1</sup>, Yuanfang Chen<sup>1,2,\*</sup>, Qiang Lin<sup>4</sup>, Kun Yang<sup>2,3</sup>, Gyu Myoung Lee<sup>5</sup>

<sup>1</sup>School of Cyberspace, Hangzhou Dianzi University, Hangzhou 310018, China

<sup>2</sup>The State Key Laboratory of Blockchain and Data Security, Zhejiang University, Hangzhou 310058, China

<sup>3</sup>College of Computer Science and Technology, Zhejiang University, Hangzhou 310058, China

<sup>4</sup>School of Computer Science and Technology, University of Science and Technology of China, Hefei 230026, China

<sup>5</sup>School of Computer Science and Mathematics, Liverpool John Moores University, Liverpool, L3 3AF, UK

\*Corresponding author: Yuanfang Chen. Email: yuanfang.chen.tina@gmail.com

1 **ABSTRACT:** The deployment of supervised anomaly detection is typically limited by the high cost of  
2 annotation, privacy constraints, and the scarcity of anomalous samples. These constraints have motivated the  
3 use of vision-language pre-trained models for zero-shot anomaly detection. However, existing CLIP-based  
4 methods still face three limitations: a shared set of prompts is applied across feature layers, anomaly maps  
5 are fused by fixed strategies, and image-level anomaly scores are determined solely by global image-text  
6 similarity. These limitations reduce the accuracy of pixel-level localization and weaken the reliability of  
7 image-level anomaly prediction. To overcome these limitations, LaRP-CLIP is proposed. It introduces  
8 layer-aware prompt decoupling to better match feature layers with different semantic characteristics,  
9 adaptive fusion with error-prior-guided local refinement to produce cleaner and more precise anomaly  
10 maps, and a prototype branch to improve image-level scoring. Experiments on four industrial datasets and  
11 seven medical datasets show that LaRP-CLIP achieves strong performance in both image-level detection and  
12 pixel-level localization.

13 **KEYWORDS:** Zero-shot anomaly detection, Vision-language models, Layer-aware prompts, Local  
14 refinement, Prototype branch.

---

## 15 1 Introduction

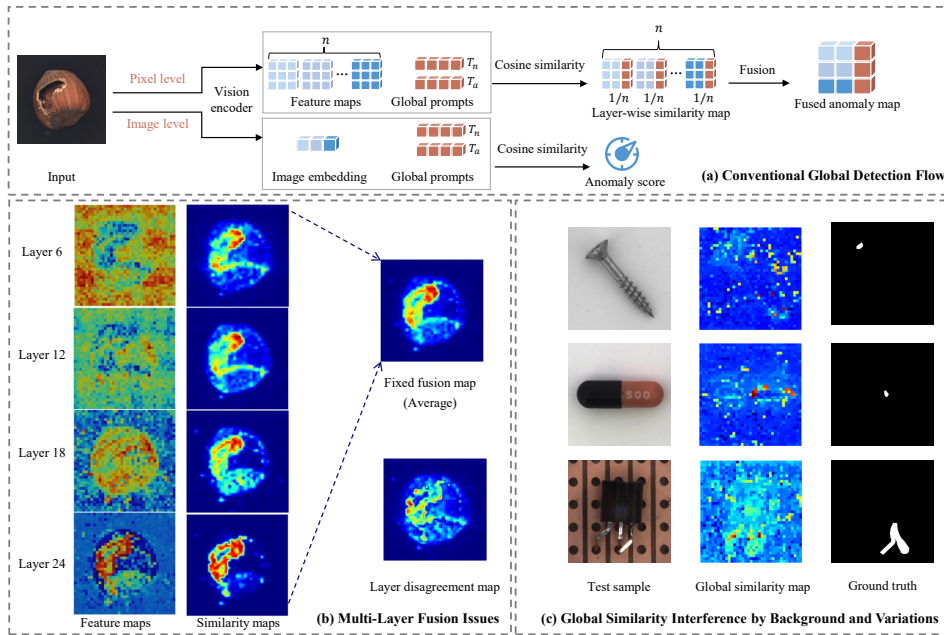
16 Anomaly detection aims to identify samples that deviate from normal patterns and plays  
17 an important role in industrial quality inspection, medical image analysis, and intelligent  
18 manufacturing [1–5]. Traditional supervised anomaly detection methods usually require a large  
19 number of normal and anomalous samples for each new product category during training [6–8]. In  
20 real industrial scenarios, however, data collection is costly, privacy constraints are common, and  
21 anomalous samples are often extremely scarce [9,10]. Therefore, this paper focuses on zero-shot  
22 anomaly detection, where the goal is to detect anomalies in unseen target domains without requiring  
23 target-domain training samples [11–15].

24 To reduce dependence on target-domain annotation, zero-shot anomaly detection (ZSAD)  
25 has attracted increasing attention in recent years [16,17]. This paradigm relies on pre-trained  
26 vision-language models, such as Contrastive Language-Image Pre-training (CLIP) [18], and

27 performs anomaly detection without requiring anomalous samples from the target domain. Among  
 28 existing ZSAD approaches, CLIP-based methods have drawn particular interest because of their  
 29 strong cross-modal representation capability [19–24]. As shown in Figure 1(a), these methods  
 30 usually extract multi-layer features from a ViT backbone, compute similarity scores with a shared  
 31 set of text prompts to obtain layer-wise anomaly maps, fuse these maps for pixel-level prediction,  
 32 and use global image features for image-level anomaly scoring.

33 Despite recent progress, existing CLIP-based ZSAD methods still exhibit three major limitations.  
 34 At the pixel level, all feature layers typically share the same set of text prompts, as illustrated  
 35 in Figure 1(b). This design is suboptimal because shallow layers mainly capture texture details,  
 36 whereas deep layers encode higher-level structure and semantics. Using the same prompts across  
 37 all layers therefore introduces a mismatch between prompt semantics and feature characteristics,  
 38 which degrades the quality of the resulting anomaly maps. Moreover, layer-wise anomaly maps  
 39 are usually combined by a fixed fusion strategy that cannot adapt to the content of each input  
 40 image. Responses from different layers may also conflict on the same defect, resulting in blurred  
 41 boundaries and increased noise in the fused maps.

42 At the image level, existing methods usually rely only on global image–text similarity, as  
 43 shown in Figure 1(c). When the input image contains a complex background or normal samples  
 44 exhibit slight appearance variations, the global similarity score can be easily affected by irrelevant  
 45 regions or benign changes. This weakens the robustness of image-level anomaly scoring.



**Figure 1:** Typical pipeline and three major limitations of existing CLIP-based zero-shot anomaly detection methods. (a) Overall workflow. (b) Pixel-level limitations: all layers share the same text prompts and use fixed fusion. (c) Image-level limitation: reliance on global image–text similarity.

46 To address these limitations, LaRP-CLIP is proposed. The main contributions of this paper are  
 47 summarized as follows:

- 48 • **Layer-Aware Prompt Decoupling:** The learnable prompts are divided into three independent  
49 groups, namely Global Prompts, Shallow Prompts, and Deep Prompts. This design allows  
50 different feature layers to interact with prompts that better match their semantic characteristics,  
51 thereby reducing the mismatch between shallow-layer texture information and deep-layer  
52 semantic representations.
- 53 • **Adaptive Weighted Fusion and Local Refinement:** Learnable layer-wise weights are  
54 introduced to adaptively fuse the anomaly maps from multiple layers. In addition, an  
55 Error-Prior Guided Local Relation Refiner is designed to perform local correction by modeling  
56 inter-layer response discrepancies, producing anomaly maps with clearer boundaries and less  
57 noise.
- 58 • **Prototype Branch:** During inference, high-confidence normal patch features are selected from  
59 the current test image to construct an image-specific normal prototype. This prototype is  
60 then combined with the global similarity score, which helps reduce the influence of complex  
61 backgrounds and normal appearance variations and improves the robustness of image-level  
62 anomaly scoring.

63 The remainder of this paper is organized as follows. Section 2 reviews related work. Section 3  
64 presents the proposed method. Section 4 reports the experimental results and analysis. Section 5  
65 concludes the paper.

## 66 2 Related Work

### 67 2.1 Traditional Anomaly Detection Methods

68 Traditional anomaly detection methods are generally divided into supervised and  
69 unsupervised paradigms. Supervised methods rely on a sufficient number of normal and  
70 anomalous samples for each category during training [3,25]. In industrial practice, however,  
71 anomalous samples are often scarce, annotation is expensive, and data sharing may be restricted by  
72 privacy or production constraints. These factors greatly hinder the rapid deployment of anomaly  
73 detection models on new product lines.

74 Unsupervised methods are usually trained with normal samples only and identify anomalies  
75 through reconstruction error [26], autoencoding mechanisms [27], or one-class modeling [28].  
76 Although such methods avoid the need for anomalous training data, they are still tied to the  
77 distribution of the training domain. Once the product type, surface texture, or imaging condition  
78 changes, their performance often degrades noticeably. This weak cross-domain generalization  
79 makes them less suitable for settings that require rapid deployment on unseen categories.

### 80 2.2 Zero-Shot Anomaly Detection

81 To reduce dependence on target-domain annotation, zero-shot anomaly detection (ZSAD) has  
82 received increasing attention in recent years [29,30]. The goal of ZSAD is to detect anomalies in a  
83 target domain without using anomalous samples from that domain during training [16]. Earlier  
84 studies mainly relied on reconstruction-based schemes or distribution modeling. In complex  
85 industrial scenes, however, reconstruction quality is easily affected by normal appearance variation,  
86 which weakens the distinction between normal samples and true anomalies [17,18].

87 With the development of large-scale foundation models, recent ZSAD research has shifted  
88 toward transfer-based solutions. Instead of learning anomaly-specific representations from

**Table 1:** Comparison of representative CLIP-based zero-shot anomaly detection methods.

Method	Venue	Prompt Design	Layer-specific Prompting	Adaptive Fusion	Local Refinement	Prototype Guidance
AnomalyCLIP [20]	ICLR'24	Learnable object-agnostic	✗	✗	✗	✗
WinCLIP [19]	CVPR'23	Hand-crafted templates	✗	✗	✗	✗
AA-CLIP [24]	CVPR'25	Anomaly-aware prompts	✗	✗	✗	✗
GenCLIP [21]	PR'26	Generalized prompts	✗	✗	✗	✗
<b>LaRP-CLIP</b>	–	Layer-aware prompts	✓	✓	✓	✓

89 target-domain data, these methods attempt to transfer general visual or cross-modal knowledge  
 90 acquired during large-scale pre-training to downstream anomaly detection tasks. This shift has  
 91 made vision-language models an important direction for zero-shot anomaly detection.

92 ZSAD is related to zero-shot learning (ZSL), where semantic knowledge is transferred to  
 93 recognize unseen classes. Recent ZSL studies, such as SVIP [31], further improve visual-semantic  
 94 alignment by modeling semantically contextualized visual patches. However, unlike ZSL  
 95 that focuses on unseen category recognition, ZSAD aims to detect category-agnostic abnormal  
 96 patterns in unseen domains. LaRP-CLIP adapts this semantic transfer principle to anomaly  
 97 detection through normal/abnormal cross-modal alignment, layer-aware prompt decoupling,  
 98 and image-specific prototype construction.

### 99 2.3 CLIP-based Zero-Shot Anomaly Detection Methods

100 The emergence of CLIP [18] has accelerated the development of CLIP-based ZSAD methods  
 101 because of its strong image–text alignment capability. WinCLIP [19] is an early representative  
 102 method that performs zero-shot anomaly classification and segmentation through a window-based  
 103 inference scheme and prompt-template ensemble. AnomalyCLIP [20] further introduces  
 104 object-agnostic learnable prompts and optimizes them with both image-level and pixel-level  
 105 supervision, showing that CLIP can be adapted more effectively to zero-shot anomaly detection.

106 Subsequent studies have mainly focused on improving prompt design, feature utilization,  
 107 and anomaly discrimination. AA-CLIP [24] strengthens anomaly awareness in both the textual  
 108 and visual spaces, thereby improving the discrimination between normal and abnormal patterns.  
 109 GenCLIP [21] further explores prompt generalization and multi-layer feature usage to improve  
 110 transferability across datasets. Despite this progress, most existing methods still rely on a single  
 111 prompt form across different layers and use fixed fusion strategies, which limits their ability to  
 112 handle layer-wise semantic differences and image-specific normal variations.

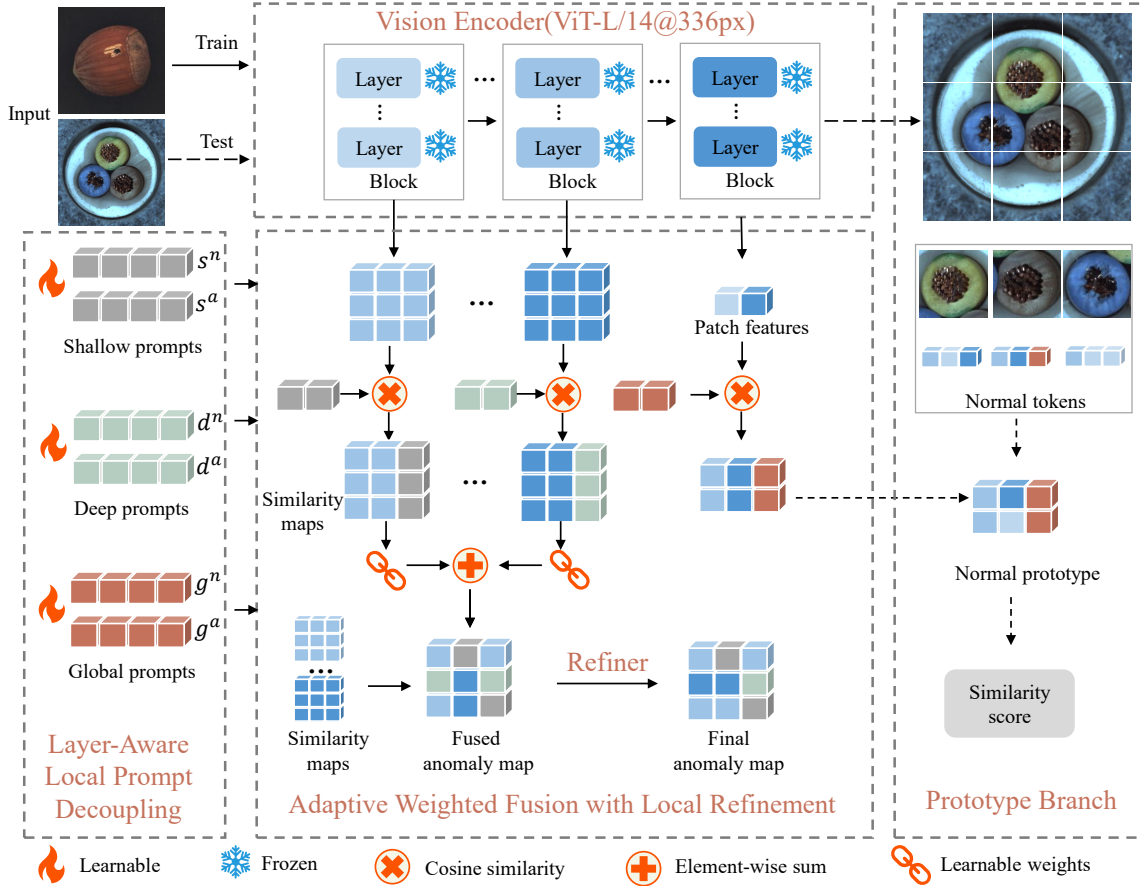
113 Table 1 summarizes representative CLIP-based ZSAD methods and highlights the differences  
 114 from LaRP-CLIP.

115 Decoupled prompts and local refinement have been explored in prompt learning and  
 116 segmentation literature but are not claimed as isolated contributions. The novelty of LaRP-CLIP  
 117 lies in adapting these ideas for CLIP-based ZSAD: prompt decoupling is driven by the semantic  
 118 hierarchy of ViT layers, local refinement is guided by inter-layer anomaly-response disagreement,  
 119 and image-level scoring is enhanced by an inference-time normal prototype. This approach directly  
 120 addresses cross-layer semantic mismatch, fixed map fusion, and unreliable global-only scoring in  
 121 existing CLIP-based ZSAD methods.

## 122 3 Method

### 123 3.1 Overall Architecture

124 LaRP-CLIP is designed for both pixel-level anomaly localization and image-level anomaly  
 125 classification in the zero-shot setting. As shown in Figure 2, given an input image  $x \in \mathbb{R}^{3 \times H \times W}$ ,  
 126 the visual encoder extracts a global image feature together with multi-layer local patch features.  
 127 Meanwhile, the prompt learner generates text prompts for different semantic levels, which are  
 128 then projected into a unified semantic space by the text encoder. The method contains two  
 129 complementary branches: a pixel-level branch for anomaly localization and an image-level branch  
 130 for robust anomaly scoring.



**Figure 2:** LaRP-CLIP consists of three main components: layer-aware prompt decoupling (left), adaptive weighted fusion with local refinement (middle), and the Prototype Branch for robust image-level scoring (right).

131 Let the global feature extracted by the visual encoder be denoted by  $f_g \in \mathbb{R}^D$ , and let the  
 132 selected multi-layer patch features be

$$F = \{F^1, F^2, \dots, F^L\}, \quad F^l \in \mathbb{R}^{N_l \times D}, \quad (1)$$

133 where  $N_l$  is the number of patches at the  $l$ -th layer and  $D$  is the feature dimension.

134 The pixel-level branch is built upon three steps: layer-aware semantic alignment, adaptive  
 135 cross-layer fusion, and local refinement. To better match the heterogeneous semantics of different  
 136 ViT layers, the local prompts are decoupled into shallow prompts and deep prompts, allowing  
 137 patch features from different layers to be aligned with anomaly semantics at suitable granularities.  
 138 Based on the resulting layer-wise anomaly responses, learnable layer weights are introduced to  
 139 perform adaptive fusion across layers. A lightweight Error-Prior Guided Local Relation Refiner is  
 140 then applied to the fused anomaly map to suppress noisy responses and improve local boundary  
 141 quality.

142 The image-level branch is introduced to improve the robustness of anomaly scoring. During  
 143 inference, a Prototype Branch constructs an image-specific normal prototype from high-confidence  
 144 normal patches selected from the current test image. This prototype serves as an intra-image  
 145 normal reference and is combined with the global similarity score, so that image-level anomaly  
 146 prediction becomes less sensitive to background clutter and benign appearance variation.

### 147 3.2 Layer-Aware Local Prompt Learning

148 Local anomaly representations exhibit clear hierarchical differences across ViT layers. Shallow  
 149 features mainly encode edges, textures, and other high-frequency local patterns, and are therefore  
 150 more suitable for describing fine-grained anomalies such as scratches, cracks, and stains. By  
 151 contrast, deep features place greater emphasis on structure, shape, and semantic consistency,  
 152 making them more suitable for structural anomalies such as missing parts, misalignment, and  
 153 deformation. If all layers share a single set of local prompts, features from different semantic levels  
 154 are forced into the same text space, which leads to a clear semantic granularity mismatch.

155 To alleviate this issue, LaRP-CLIP divides the learnable prompts into three groups, namely  
 156 global prompts, shallow prompts, and deep prompts. After the text encoder, the corresponding  
 157 text embeddings are written as

$$158 [(\hat{g}^n, \hat{g}^a), (\hat{s}^n, \hat{s}^a), (\hat{d}^n, \hat{d}^a)] \in \mathbb{R}^{2 \times D}, \quad (2)$$

159 where the superscripts  $n$  and  $a$  denote the normal and abnormal classes, respectively.

160 The layers used for local anomaly modeling are divided into a shallow set  $S$  and a deep set  $D$ ,  
 161 satisfying

$$162 S \cup D = \{1, 2, \dots, L\}, \quad S \cap D = \emptyset. \quad (3)$$

163 For the patch features  $F^l$  at the  $l$ -th layer, each patch feature is first normalized by

$$164 \hat{f}_i^l = \frac{f_i^l}{\|f_i^l\|_2}, \quad (4)$$

165 where  $f_i^l$  denotes the  $i$ -th patch feature at layer  $l$ . The text embeddings are normalized in the  
 166 same way. If  $l \in S$ , the shallow prompt embeddings  $(\hat{s}^n, \hat{s}^a)$  are used; if  $l \in D$ , the deep prompt  
 167 embeddings  $(\hat{d}^n, \hat{d}^a)$  are used.

165 The similarity between the  $i$ -th patch and the corresponding normal or abnormal text  
166 embedding is defined as

$$z_{i,c}^l = \tau (\hat{f}_i^l)^\top \hat{e}_c^l, \quad c \in \{n, a\}, \quad (5)$$

167 where  $\tau$  is the temperature coefficient and  $\hat{e}_c^l$  denotes the text embedding assigned to layer  $l$ . The  
168 anomaly response of the  $i$ -th patch is then obtained by binary softmax:

$$p_{i,a}^l = \frac{\exp(z_{i,a}^l)}{\exp(z_{i,n}^l) + \exp(z_{i,a}^l)}. \quad (6)$$

169 The anomaly map of the  $l$ -th layer is written as

$$M^l = \text{Reshape}(\{p_{i,a}^l\}_{i=1}^{N_l}). \quad (7)$$

170 For image-level scoring, the global feature interacts only with the text embeddings  
171 corresponding to the global prompts. The resulting text-driven anomaly score is defined as  
172

$$s_{\text{text}} = \frac{\exp(\tau_g \hat{f}_g^\top \hat{g}^a)}{\exp(\tau_g \hat{f}_g^\top \hat{g}^n) + \exp(\tau_g \hat{f}_g^\top \hat{g}^a)}. \quad (8)$$

173 With this design, shallow patch features are aligned with shallow anomaly semantics,  
174 deep patch features are aligned with deep anomaly semantics, and the global feature is used  
175 for image-level anomaly judgment. In this way, the cross-modal alignment process becomes  
176 layer-aware, which helps reduce the semantic mismatch introduced by using a shared prompt  
177 across all feature layers.

### 178 *3.3 Layer-Adaptive Fusion and Local Relation Refinement*

179 The layer-wise anomaly maps  $\{M^l\}_{l=1}^L$  differ in detail sensitivity, structural stability, and noise  
180 level. Fixed fusion or equal-weight averaging cannot fully exploit the complementary information  
181 carried by different layers. To address this issue, LaRP-CLIP introduces a set of learnable fusion  
182 parameters  $\{\alpha^l\}_{l=1}^L$ , which are normalized by softmax to obtain the layer weights:

$$\omega_l = \frac{\exp(\alpha_l)}{\sum_{k=1}^L \exp(\alpha_k)}, \quad l = 1, \dots, L. \quad (9)$$

183 The fused anomaly map is then computed as

$$M_{\text{fused}} = \sum_{l=1}^L \omega_l M^l. \quad (10)$$

184 This design allows the contribution of each layer to be learned from data, rather than being  
185 fixed in advance.

186 Although adaptive fusion improves the integration of multi-layer responses, the fused map  
187 may still contain local noise or imprecise boundaries. This issue becomes more evident when

188 different layers produce inconsistent responses in the same region. To capture such disagreement,  
 189 an error prior is introduced to measure the deviation of each layer from the fused result. For the  
 190  $l$ -th layer, the deviation map is defined as

$$E^l = |M^l - M_{\text{fused}}|. \quad (11)$$

191 The aggregated error prior is then written as

$$E_{\text{prior}} = \sum_{l=1}^L \omega_l E^l. \quad (12)$$

192 The error prior provides a local measure of cross-layer inconsistency. Regions with similar  
 193 responses across layers tend to have small values, whereas regions with conflicting responses yield  
 194 larger values. Based on this cue, LaRP-CLIP concatenates the fused anomaly map with the error  
 195 prior and feeds them into a lightweight local relation refinement module. This module adopts a  
 196  $3 \times 3$  convolutional structure to predict a residual correction:

$$\Delta M = \psi([M_{\text{fused}}, E_{\text{prior}}]). \quad (13)$$

197 The final pixel-level anomaly map is given by

$$M_{\text{final}} = M_{\text{fused}} + \Delta M. \quad (14)$$

198 The refinement module is used to correct local errors in the fused anomaly map rather than to  
 199 regenerate the entire prediction. It mainly adjusts regions with blurred boundaries, discontinuous  
 200 responses, or noticeable local noise, while preserving the main structure provided by the preceding  
 201 layer-wise fusion stage.

### 202 **3.4 Prototype Branch**

203 Relying only on global text matching is often insufficient for stable image-level anomaly  
 204 scoring, because it does not fully reflect the local distribution of anomaly evidence within an image.  
 205 This limitation becomes more apparent when anomalous regions are small, the background is  
 206 complex, or normal samples exhibit noticeable appearance variation. To improve the robustness of  
 207 image-level scoring, a Prototype Branch is introduced during inference. This branch constructs a  
 208 normal prototype from high-confidence normal regions in the current test image and uses it as an  
 209 additional reference for anomaly judgment.

210 The layer-wise anomaly maps are first accumulated to obtain a text-driven anomaly map:

$$M_{\text{text}} = \sum_{l=1}^L M^l. \quad (15)$$

211 The patch features from the last layer, denoted by  $\{p_i\}_{i=1}^N$ , are then used as local semantic  
 212 representations. The map  $M_{\text{text}}$  is downsampled to the patch grid to produce an anomaly score  $q_i$

213 for each patch. Based on these scores, a set of high-confidence normal patches is selected, and its  
 214 index set is denoted by  $\mathcal{N}$ . The normal prototype is computed as

$$p_{\text{norm}} = \frac{\sum_{i \in \mathcal{N}} \omega_i p_i}{\sum_{i \in \mathcal{N}} \omega_i}, \quad (16)$$

215 where  $\omega_i$  is derived from  $q_i$ , so that patches with lower anomaly scores receive larger weights.

216 The deviation of each patch from the normal prototype is then measured by

$$\mu_i = 1 - \cos(p_i, p_{\text{norm}}). \quad (17)$$

217 These deviations are reshaped into a two-dimensional map to form the prototype-guided anomaly  
 218 map, which is further combined with the text-driven anomaly map to obtain an auxiliary anomaly  
 219 map for image-level scoring. Since image-level prediction is more influenced by the most salient  
 220 anomalous regions than by the average response over the whole image, top- $k$  pooling is applied to  
 221 the fused map to obtain the auxiliary score  $s_{\text{map}}$ . The final image-level anomaly score is defined as

$$s_{\text{final}} = \frac{s_{\text{text}} + \lambda s_{\text{map}}}{1 + \lambda}, \quad (18)$$

222 where  $\lambda$  is a balancing coefficient.

223 Unlike a fixed normal template shared across all images, the Prototype Branch derives a normal  
 224 reference directly from the current test image. This makes the image-level score less sensitive to  
 225 background clutter and benign appearance variation. The branch is used only for image-level  
 226 scoring and does not affect the generation of pixel-level anomaly maps.

227 The Prototype Branch is excluded from training because it is designed as a non-parametric  
 228 test-time adaptation module. It constructs an image-specific normal prototype from high-confidence  
 229 normal patches in each test image, rather than learning a fixed prototype from auxiliary training  
 230 data.

### 231 3.5 Training Objective

232 To jointly optimize image-level anomaly discrimination and pixel-level anomaly localization,  
 233 LaRP-CLIP uses a combined training objective. The overall loss is defined as

$$\mathcal{L} = \mathcal{L}_{\text{img}} + \mathcal{L}_{\text{pixel}}, \quad (19)$$

234 where  $\mathcal{L}_{\text{img}}$  is the image-level loss and  $\mathcal{L}_{\text{pixel}}$  is the pixel-level loss.

235 For image-level supervision, let  $y \in \{0, 1\}$  denote the image label and let  $s_{\text{text}}$  denote the global  
 236 text-driven anomaly score. The image-level loss is defined by binary cross-entropy:

$$\mathcal{L}_{\text{img}} = -y \log s_{\text{text}} - (1 - y) \log(1 - s_{\text{text}}). \quad (20)$$

237 For pixel-level supervision, the intermediate outputs are not treated as independent prediction  
 238 targets. Instead, the loss is divided according to the two stages of the pixel-level branch, namely  
 239 alignment and refinement:

$$\mathcal{L}_{\text{pixel}} = \mathcal{L}_{\text{align}} + \lambda_r \mathcal{L}_{\text{ref}}, \quad (21)$$

240 where  $\mathcal{L}_{\text{align}}$  supervises the front-end layer-aware semantic alignment and adaptive fusion,  $\mathcal{L}_{\text{ref}}$   
 241 supervises the refinement stage, and  $\lambda_r$  is the weight assigned to the refinement term.

242 The alignment loss is applied to the shallow-layer anomaly maps, the deep-layer anomaly  
 243 maps, and the fused anomaly map:

$$\begin{aligned} \mathcal{L}_{\text{align}} = & \frac{1}{|S|} \sum_{l \in S} \mathcal{L}_{\text{seg}}(M^l, Y) \\ & + \frac{1}{|D|} \sum_{l \in D} \mathcal{L}_{\text{seg}}(M^l, Y) \\ & + \lambda_f \mathcal{L}_{\text{seg}}(M_{\text{fused}}, Y), \end{aligned} \quad (22)$$

244 where  $Y$  is the pixel-level ground-truth mask,  $M^l$  is the anomaly map at layer  $l$ ,  $M_{\text{fused}}$  is the  
 245 adaptively fused anomaly map, and  $S$  and  $D$  denote the shallow-layer set and deep-layer set,  
 246 respectively. The function  $\mathcal{L}_{\text{seg}}(\cdot)$  denotes the segmentation loss, implemented as a combination  
 247 of focal loss and dual-channel Dice loss for the normal and abnormal classes. The coefficient  $\lambda_f$   
 248 controls the relative contribution of the fused-map supervision.

249 This design serves two purposes. First, separate supervision on shallow and deep anomaly  
 250 maps encourages the decoupled prompts to learn anomaly semantics at different granularities.  
 251 Second, direct supervision on  $M_{\text{fused}}$  ensures that the fusion stage itself produces a discriminative  
 252 anomaly map before refinement.

253 The refinement loss is applied only to the final anomaly map:

$$\mathcal{L}_{\text{ref}} = \mathcal{L}_{\text{seg}}(M_{\text{final}}, Y) + \lambda_{\Delta} \|M_{\text{final}} - M_{\text{fused}}\|_1, \quad (23)$$

254 where  $\lambda_{\Delta}$  is the weight of the residual regularization term. This term limits the difference between  
 255 the refined anomaly map and the fused anomaly map, thereby constraining the magnitude of the  
 256 correction introduced by the refinement module.

257 The Prototype Branch is used only during inference and is therefore not included in the training  
 258 objective.

### 259 *3.6 Training and Inference Procedure*

260 To improve the readability and reproducibility of the proposed method, the pseudo code of  
 261 LaRP-CLIP is provided in Algorithm 1. The algorithm summarizes the complete training and  
 262 inference procedure, including layer-aware prompt matching, adaptive fusion with local refinement,  
 263 and the Prototype Branch for image-level scoring. During training, the loss in Eq. (19) is used to  
 264 optimize the learnable parameters. During inference, all learnable parameters are fixed, and only  
 265 the forward pass is performed to obtain the final anomaly scores.

---

**Algorithm 1** The pseudo code of LaRP-CLIP for zero-shot anomaly detection.

---

**Require:** Input image  $x \in \mathbb{R}^{3 \times H \times W}$ , pre-trained visual and text encoders, learnable prompts  $\{\alpha^l, g^n, g^a, s^n, s^a, d^n, d^a\}$ , refinement module  $\psi$ , balancing coefficient  $\lambda$

**Ensure:** Pixel-level anomaly map  $M_{\text{final}}$ , image-level anomaly score  $s_{\text{final}}$

- 1: Extract global feature  $f_g$  and multi-layer patch features  $F = \{F^l\}_{l=1}^L$  {Eq. (1)}
- 2: Generate and normalize global, shallow, and deep text embeddings:  $(\hat{g}^n, \hat{g}^a), (\hat{s}^n, \hat{s}^a), (\hat{d}^n, \hat{d}^a)$  {Eq. (2)}
- 3: Partition the selected layers into shallow set  $S$  and deep set  $D$  {Eq. (3)}
- 4: **for** each layer  $l = 1, \dots, L$  **do**
- 5: Normalize patch features:  $\hat{f}_i^l = f_i^l / \|f_i^l\|_2$  {Eq. (4)}
- 6: Select the corresponding prompt embedding according to  $l \in S$  or  $l \in D$
- 7: Compute similarity scores:  $z_{i,c}^l = \tau(\hat{f}_i^l)^\top \hat{e}_c^l$  {Eq. (5)}
- 8: Compute the anomaly response  $p_{i,a}^l$  and form the anomaly map  $M^l$  {Eqs. (6), (7)}
- 9: **end for**
- 10: Compute normalized layer weights  $\omega_l$  by softmax {Eq. (9)}
- 11: Fuse the layer-wise anomaly maps:  $M_{\text{fused}} = \sum_{l=1}^L \omega_l M^l$  {Eq. (10)}
- 12: Compute deviation maps and error prior:  $E^l = |M^l - M_{\text{fused}}|, E_{\text{prior}} = \sum_{l=1}^L \omega_l E^l$  {Eqs. (11), (12)}
- 13: Predict residual correction and obtain the final anomaly map:  $\Delta M = \psi([M_{\text{fused}}, E_{\text{prior}}]), M_{\text{final}} = M_{\text{fused}} + \Delta M$  {Eqs. (13), (14)}
- 14: Compute the global text-driven anomaly score  $s_{\text{text}}$  {Eq. (8)}
- 15: Accumulate the layer-wise anomaly maps:  $M_{\text{text}} = \sum_{l=1}^L M^l$  {Eq. (15)}
- 16: Select high-confidence normal patches and compute the normal prototype  $p_{\text{norm}}$  {Eq. (16)}
- 17: Compute the prototype-guided auxiliary score  $s_{\text{map}}$  by top- $k$  pooling
- 18: Compute the final image-level score:  $s_{\text{final}} = (1 - \lambda)s_{\text{text}} + \lambda s_{\text{map}}$  {Eq. (18)}
- 19: **if** training **then**
- 20: Compute the total loss  $\mathcal{L} = \mathcal{L}_{\text{img}} + \mathcal{L}_{\text{pixel}}$  {Eqs. (19)–(23)}
- 21: **end if**
- 22: **return**  $M_{\text{final}}, s_{\text{final}}$

---

## 266 4 Experiments

### 267 4.1 Experimental Settings

268 The experiments are evaluated on the MVTEC AD [32], VisA [23], MPDD [33], and BTAD  
 269 [34] industrial defect datasets, as well as the HeadCT [35], BrainMRI [35], Br35H [36], ISIC [37],  
 270 CVC-ColonDB [38], CVC-ClinicDB [38], and TN3K [39] medical imaging datasets. For evaluation  
 271 metrics, pixel-level performance is measured by pixel-level AUROC and PRO, while image-level  
 272 performance is measured by image-level AUROC and AP.

273 All experiments were conducted on a single NVIDIA RTX 4090 GPU (24 GB VRAM) using  
 274 Python 3.12.7 and the PyTorch 2.0.0 framework. The Adam optimizer was employed with an initial  
 275 learning rate of  $1 \times 10^{-3}$ , a batch size of 8, and a total of 15 training epochs. To ensure reproducibility,  
 276 the random seed was fixed at 111. The visual backbone network adopts the pre-trained CLIP model  
 277 (ViT-L/14@336px), with the input resolution uniformly set to  $518 \times 518$ . During training, only the  
 278 parameters of the learnable prompt learner are optimized, while the CLIP visual encoder remains  
 279 frozen.

**Table 2:** Quantitative comparison on industrial defect datasets (MVTec AD, VisA, MPDD, BTAD).

Methods	CLIP [18]	CLIP-AC [18]	WinCLIP [19]	VAND [22]	CoOp [40]	AnomalyCLIP [20]	AACLIP [24]	GenCLIP [21]	LaRP-CLIP
<b>Image-level (AUROC, AP)↑</b>									
MVTec AD	(74.1, 87.6)	(71.5, 86.4)	(91.8, <b>96.5</b> )	(86.1, 93.5)	(88.8, 94.8)	<b>(91.5, 96.2)</b>	(90.5, 94.9)	(90.9, 96.1)	<b>(92.8, 96.5)</b>
VisA	(66.4, 71.5)	(65.0, 70.1)	(78.1, 81.2)	(78.0, 81.4)	(62.8, 68.1)	<b>(82.1, 85.4)</b>	(79.0, 81.9)	(83.3, 87.5)	<b>(86.0, 88.0)</b>
MPDD	(54.3, 65.4)	(56.2, 66.0)	(63.6, 69.9)	(73.0, 80.2)	(55.1, 64.2)	<b>(77.0, 82.0)</b>	(56.0, 66.6)	(73.7, 79.6)	<b>(77.7, 80.7)</b>
BTAD	(34.5, 52.5)	(51.0, 62.1)	(68.2, 70.9)	(73.6, 68.6)	(66.8, 77.4)	(88.3, 87.3)	<b>(92.9, 96.5)</b>	(90.0, <b>96.9</b> )	<b>(91.6, 94.3)</b>
<b>Pixel-level (AUROC, PRO)↑</b>									
MVTec AD	(38.4, 11.3)	(38.2, 11.6)	(85.2, 64.6)	(87.6, 44.0)	(33.3, 6.7)	(91.1, 81.4)	<b>(92.3, 85.7)</b>	<b>(92.7, 88.1)</b>	(91.5, <b>85.8</b> )
VisA	(46.6, 14.8)	(47.8, 17.3)	(79.6, 56.8)	(94.2, 86.8)	(24.2, 3.8)	<b>(95.5, 87.0)</b>	(94.8, 83.7)	<b>(92.7, 88.1)</b>	<b>(96.7, 91.7)</b>
MPDD	(62.1, 33.0)	(58.7, 29.1)	(76.4, 48.9)	(94.1, 83.2)	(15.4, 2.3)	<b>(96.5, 88.7)</b>	(95.7, 84.8)	(96.2, <b>89.3</b> )	<b>(97.6, 90.1)</b>
BTAD	(30.6, 4.4)	(32.8, 8.3)	(72.7, 27.3)	(60.8, 25.0)	(28.6, 3.8)	(94.2, 74.8)	<b>(94.5, 79.1)</b>	(93.6, 75.6)	<b>(94.8, 77.5)</b>

**Table 3:** Quantitative comparison on medical imaging datasets (HeadCT, BrainMRI, Br35H, ISIC, CVC-ColonDB, CVC-ClinicDB, TN3K).

Methods	CLIP [18]	CLIP-AC [18]	WinCLIP [19]	VAND [22]	CoOp [40]	AnomalyCLIP [20]	AACLIP [24]	GenCLIP [21]	LaRP-CLIP
<b>Image-level (AUROC, AP)↑</b>									
HeadCT	(56.5, 58.4)	(60.0, 60.7)	(81.8, 80.2)	(89.1, 89.4)	(78.4, 78.8)	(93.4, 91.6)	(84.8, 87.4)	<b>(96.1, 96.2)</b>	<b>(96.7, 97.3)</b>
BrainMRI	(73.9, 81.7)	(80.6, 86.4)	(86.6, 91.5)	(89.3, 90.9)	(61.3, 44.9)	(90.3, 92.2)	(90.3, 93.3)	<b>(92.4, 94.4)</b>	<b>(95.6, 96.4)</b>
Br35H	(78.4, 78.8)	(82.7, 81.3)	(80.5, 82.2)	(93.1, 92.9)	(86.0, 87.5)	(94.6, 94.7)	(86.9, 88.6)	<b>(95.0, 95.0)</b>	<b>(97.7, 97.6)</b>
<b>Pixel-level (AUROC, PRO)↑</b>									
ISIC	(33.1, 5.8)	(36.0, 7.7)	(83.3, 55.1)	(89.4, 77.2)	(51.7, 15.9)	(89.7, 78.4)	<b>(94.0, 79.7)</b>	<b>(93.1, 88.1)</b>	(92.5, <b>83.0</b> )
CVC-ColonDB	(49.5, 15.8)	(49.5, 11.5)	(70.3, 32.5)	(78.4, 64.6)	(40.5, 2.6)	<b>(81.9, 71.3)</b>	(80.2, 56.6)	(80.3, <b>87.8</b> )	(81.8, 71.0)
CVC-ClinicDB	(47.5, 18.9)	(48.5, 12.6)	(51.2, 13.8)	(80.5, 60.7)	(34.8, 2.4)	(82.9, <b>67.8</b> )	<b>(81.3, 53.2)</b>	<b>(81.3, 82.8)</b>	<b>(85.2, 67.1)</b>
TN3K	(42.3, 7.3)	(35.6, 5.2)	(70.7, 39.8)	(73.6, 37.8)	(34.0, 9.5)	<b>(81.5, 50.4)</b>	(75.9, 44.7)	(73.8, <b>80.6</b> )	<b>(82.5, 51.9)</b>

280 The prompt learner uses a learnable prompt length of 12, a text embedding depth of 9, and  
 281 feature layers selected as [6, 12, 18, 24]. The temperature coefficient is  $\tau = 0.07$ , the balancing  
 282 coefficient is  $\lambda = 1.0$ ,  $\lambda_f = 1$ ,  $\lambda_r = 0.5$ , and  $\lambda_\Delta = 0.2$ . The Prototype Branch is enabled during the  
 283 testing phase with top- $k$  set to 8.

## 284 4.2 Main Results

285 Table 2 and Table 3 summarize the quantitative comparison between LaRP-CLIP and  
 286 state-of-the-art zero-shot anomaly detection methods on industrial defect datasets and medical  
 287 imaging datasets, respectively. The highest value of each metric (AUROC or AP/PRO) is  
 288 highlighted in **red** and the second-highest value is highlighted in **blue**.

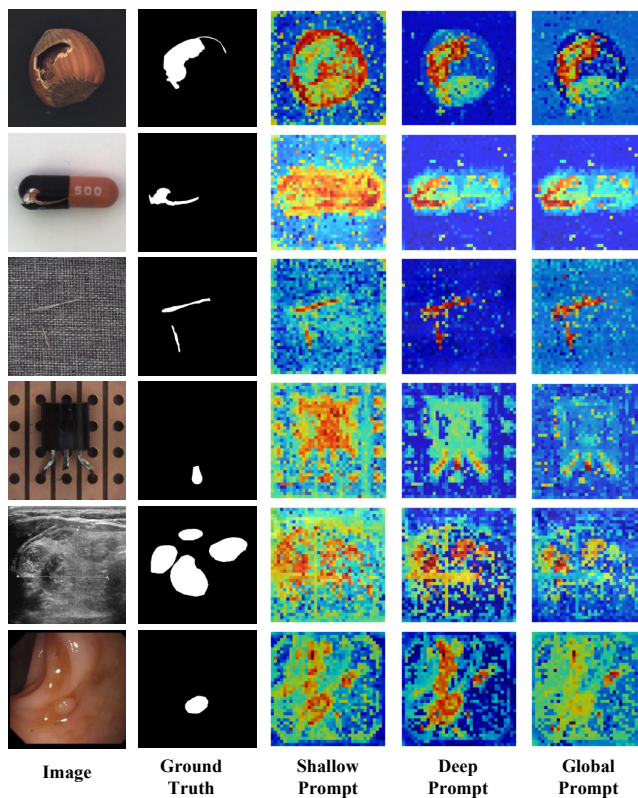
289 On the four industrial defect datasets, LaRP-CLIP achieves the best image-level performance  
 290 on MVTec AD and VisA, with AUROC/AP scores of (92.8, 96.5) and (86.0, 88.0), respectively.  
 291 On MPDD and BTAD, it also remains among the strongest competitors, yielding the second-best  
 292 image-level results. For pixel-level evaluation, LaRP-CLIP consistently ranks first or second across  
 293 all four datasets and surpasses previous CLIP-based methods in both AUROC and PRO in most  
 294 cases. In particular, compared with AnomalyCLIP, the pixel-level performance is improved from  
 295 (95.5, 87.0) to (96.7, 91.7) on VisA and from (96.5, 88.7) to (97.6, 90.1) on MPDD. The results show  
 296 that LaRP-CLIP performs consistently well across industrial datasets with diverse defect patterns  
 297 and visual characteristics.

298 On the medical imaging benchmarks, LaRP-CLIP also shows strong transferability. It achieves  
 299 the highest image-level AUROC and AP on HeadCT, BrainMRI, and Br35H, with scores of (96.7,  
 300 97.3), (95.6, 96.4), and (97.7, 97.6), respectively. For pixel-level anomaly localization, LaRP-CLIP

301 attains the best AUROC on CVC-ClinicDB and TN3K, while remaining highly competitive on ISIC  
 302 and CVC-ColonDB. It also achieves strong PRO performance, particularly on datasets where lesion  
 303 boundaries are difficult to delineate. This cross-domain performance indicates that LaRP-CLIP  
 304 remains reliable even when the target anomalies differ markedly from the training domain in both  
 305 appearance and imaging style.

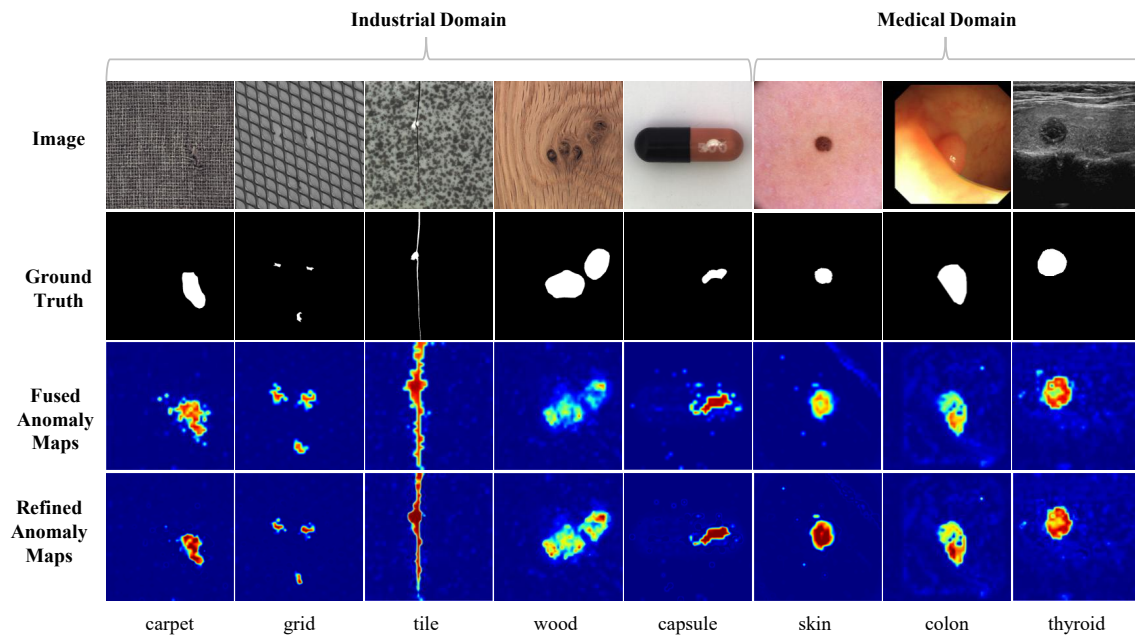
306 It is also observed that LaRP-CLIP does not achieve the best PRO on all medical datasets.  
 307 On ISIC and CVC-ColonDB, lower PRO scores indicate incomplete localization for lesions with  
 308 irregular, low-contrast, or ambiguous boundaries. This suggests that although LaRP-CLIP shows  
 309 competitive cross-domain transferability, stronger boundary-aware modeling is still needed for  
 310 challenging medical lesions.

311 **Figure 3** presents a qualitative comparison of anomaly maps produced by different prompt  
 312 strategies. The three prompts show notably different response patterns across samples. The Shallow  
 313 Prompt is more sensitive to local appearance variations and can better preserve fine-grained details,  
 314 making it effective for subtle texture anomalies. The Deep Prompt generates more compact and  
 315 semantically aligned activation on abnormal regions, which is beneficial for anomalies with  
 316 relatively clear structural semantics. In contrast, the Global Prompt frequently produces diffuse  
 317 or less discriminative responses, leading to incomplete localization or increased background  
 318 interference. This comparison indicates that different ViT layers capture anomaly cues at different  
 319 semantic granularities, while the proposed layer-aware prompt decoupling offers a more suitable  
 320 mechanism for aligning prompt representations with layer-specific features.



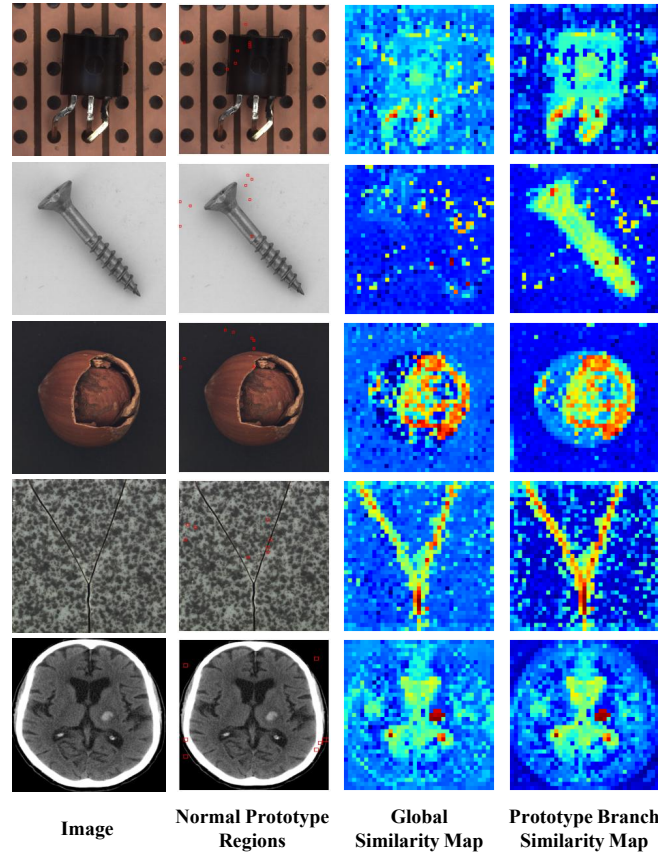
**Figure 3:** Qualitative comparison of anomaly maps generated using different prompts. From left to right: input image, ground-truth mask, anomaly map using Shallow Prompt, Deep Prompt, and Global Prompt.

321 **Figure 4** presents qualitative anomaly localization results of LaRP-CLIP on representative  
 322 samples from both industrial and medical domains. The third row shows the fused anomaly maps  
 323 obtained after multi-layer aggregation, while the fourth row shows the refined anomaly maps after  
 324 applying the Error-Prior Guided Local Relation Refiner. As observed, the fused maps can already  
 325 roughly highlight the abnormal regions, but they still contain scattered responses and imprecise  
 326 boundaries in several cases. After refinement, the anomaly regions become more compact and  
 327 better aligned with the ground-truth masks, with reduced background interference and clearer  
 328 structural outlines. This improvement is consistently visible in both industrial samples with thin or  
 329 irregular defects and medical samples with relatively compact lesion regions, indicating that the  
 330 refinement stage is beneficial for producing more precise pixel-level localization.



**Figure 4:** Qualitative anomaly localization results of LaRP-CLIP on representative samples from industrial and medical domains. From top to bottom: input image, ground-truth mask, fused anomaly map, and refined anomaly map. The refined maps exhibit clearer boundaries and reduced background interference compared with the fused maps.

331 **Figure 5** presents a qualitative comparison between the global similarity map and the  
 332 prototype-guided similarity map. The selected high-confidence normal prototype regions in  
 333 the second column are used to build an image-specific normal reference for the current test sample.  
 334 Compared with the global similarity map, which tends to exhibit scattered or distracting responses  
 335 on normal structures, the prototype-guided similarity map yields cleaner activations and better  
 336 highlights the actual abnormal regions. The improvement is visible across samples from both  
 337 industrial and medical domains, indicating that the Prototype Branch helps reduce the influence of  
 338 normal appearance variation and complex background content during anomaly localization.



**Figure 5:** Qualitative results showing the effect of the Prototype Branch. From left to right: input image, selected high-confidence normal prototype regions (red dots), global similarity map, and prototype-guided similarity map.

### 339 4.3 Ablation Study

340 To assess the contribution of each component in LaRP-CLIP, an ablation study is conducted on  
 341 MVTEC AD. The results are reported in Table 4, where image-level performance is evaluated by  
 342 AUROC and AP, and pixel-level performance is evaluated by AUROC and PRO.

343 As shown in Table 4, removing any component leads to a decline in performance, indicating  
 344 that each module contributes to the overall capability of the framework. Among all variants,  
 345 removing the Layer-Aware Prompt Decoupling causes the largest performance drop, with  
 346 image-level results decreasing from (92.8, 96.5) to (88.4, 93.7) and pixel-level results decreasing  
 347 from (91.5, 85.8) to (87.9, 80.2). This result suggests that aligning prompt representations with  
 348 the semantic granularity of different ViT layers is critical for both anomaly classification and  
 349 localization, which is also consistent with the qualitative comparison in Figure 3.

350 When the Adaptive Weighted Fusion is removed, the performance decreases to (89.7, 94.9) at  
 351 the image level and (89.1, 82.6) at the pixel level, suggesting that fixed fusion is less effective for  
 352 aggregating anomaly cues from multiple layers. Removing the Local Relation Refiner also results  
 353 in a noticeable decline, particularly in pixel-level localization, where the PRO score drops from  
 354 85.8 to 81.9. This observation is consistent with the qualitative comparison in Figure 4, where the  
 355 refinement stage produces cleaner anomaly regions and clearer boundaries.

**Table 4:** Ablation study on MVTec AD. “w/o” denotes the removal of the corresponding component.

Variant	Image-level	Pixel-level
Full LaRP-CLIP	(92.8, 96.5)	(91.5, 85.8)
w/o Layer-Aware Prompt Decoupling	(88.4, 93.7)	(87.9, 80.2)
w/o Adaptive Weighted Fusion	(89.7, 94.9)	(89.1, 82.6)
w/o Local Relation Refiner	(90.2, 95.1)	(88.4, 81.9)
w/o Prototype Branch	(90.5, 95.4)	(91.0, 85.1)

**Table 5:** Sensitivity analysis of  $\lambda$  and top- $k$  on MVTec AD. Image-level results are reported as (AUROC, AP).

Parameter	Setting	Image-level
$\lambda$	0.0	(91.6, 96.1)
$\lambda$	0.5	(92.4, 96.4)
$\lambda$	1.0	(92.8, 96.5)
$\lambda$	2.0	(92.5, 96.3)
$\lambda$	4.0	(92.1, 96.0)
top- $k$	1	(92.0, 96.1)
top- $k$	4	(92.5, 96.4)
top- $k$	8	(92.8, 96.5)
top- $k$	16	(92.6, 96.4)
top- $k$	32	(92.2, 96.1)

356 By comparison, removing the Prototype Branch leads to a relatively smaller drop, from (92.8,  
357 96.5) to (90.5, 95.4) at the image level and from (91.5, 85.8) to (91.0, 85.1) at the pixel level. Although  
358 its effect is less pronounced than that of the other modules, it still brings consistent gains, indicating  
359 that image-specific normal prototypes help stabilize anomaly scoring under appearance variation.  
360 This trend is also supported by the visual comparison in Figure 5, where the prototype-guided  
361 similarity maps show more concentrated responses on abnormal regions and less interference from  
362 normal structures.

363 The results reveal clear performance degradation when any component is removed. Removing  
364 layer-aware prompt decoupling leads to the largest drop, decreasing image-level AUROC by 4.4%  
365 and pixel-level AUROC by 3.6%, confirming its critical role in resolving semantic granularity  
366 mismatch across layers. Removing the Error-Prior Guided Local Relation Refiner causes noticeable  
367 declines in pixel-level metrics (AUROC drops by 3.1% and PRO by 4.0%), highlighting its  
368 importance in sharpening boundaries and suppressing noise. Eliminating the Prototype Branch  
369 primarily affects image-level robustness, while equal-weight fusion and the absence of error prior  
370 both result in suboptimal fusion quality.

371 To further evaluate the influence of key hyperparameters, sensitivity analyses are conducted  
372 on MVTec AD. Three factors are considered: the balancing coefficient  $\lambda$ , the top- $k$  value used in  
373 prototype-guided pooling, and the shallow/deep layer partition. Since  $\lambda$  and top- $k$  directly affect  
374 the final image-level anomaly score, their influence is evaluated using image-level AUROC and  
375 AP. The shallow/deep layer partition affects layer-aware prompt matching and anomaly map  
376 generation, and is therefore evaluated using pixel-level AUROC and PRO.

377 As shown in Table 5, LaRP-CLIP remains stable with different  $\lambda$  and top- $k$  values. When  $\lambda = 0$ ,  
378 the score relies only on global text similarity, leading to weaker detection. Increasing  $\lambda$  incorporates  
379 local anomaly evidence, improving performance, but too large a  $\lambda$  overemphasizes local responses,

**Table 6:** Sensitivity analysis of shallow/deep layer partition on MVTec AD. Pixel-level results are reported as (AUROC, PRO).

Layer Partition ( <i>S/D</i> )	Pixel-level
[6]/[12, 18, 24]	(90.7, 84.6)
[6, 12]/[18, 24]	(91.5, 85.8)
[6, 12, 18]/[24]	(91.0, 85.1)

380 slightly degrading performance. The best result is achieved at  $\lambda = 1.0$ , where global and local  
381 evidence complement each other.

382 For top- $k$ , a very small value is more sensitive to noisy local activations, while a very large  
383 value may include normal background regions and dilute anomaly evidence. The setting top- $k = 8$   
384 achieves the best performance and is therefore adopted as the default configuration.

385 Table 6 shows that assigning only layer 6 to the shallow group lacks sufficient texture  
386 supervision, while assigning layer 18 weakens the distinction between shallow and deep prompts.  
387 The best pixel-level performance is achieved with the partition [6, 12]/[18, 24], aligning with the  
388 ViT feature hierarchy: shallow layers capture texture, while deep layers encode structure and  
389 semantics.

#### 390 4.4 Time Complexity Analysis

391 Let  $N$  denote the number of image patches and  $L$  denote the number of selected feature layers.  
392 In LaRP-CLIP, the extra computation beyond the frozen CLIP backbone comes from layer-wise  
393 prompt matching, adaptive fusion, refinement, and prototype construction.

394 Layer-aware prompt matching and patch-text similarity computation across  $L$  layers require  
395  $O(LN)$ . Adaptive weighted fusion and error-prior estimation are also performed on layer-wise  
396 anomaly responses, with complexity  $O(LN)$ . The Local Relation Refiner is implemented using  
397 lightweight  $3 \times 3$  convolutions and has complexity  $O(N)$ . The Prototype Branch, including  
398 high-confidence patch selection and prototype construction, also scales as  $O(N)$ .

399 The total additional complexity is therefore  $O(LN)$ , excluding the frozen backbone. This is  
400 on the same order as existing CLIP-based methods that rely on multi-layer inference. On a single  
401 NVIDIA RTX 4090 GPU, with an input resolution of  $518 \times 518$ , LaRP-CLIP takes about 280 ms to  
402 process one image during inference.

## 403 5 Conclusion

404 This paper presents LaRP-CLIP, a zero-shot anomaly detection framework for industrial and  
405 medical images. The framework introduces layer-aware prompt decoupling, adaptive weighted  
406 fusion with an Error-Prior Guided Local Relation Refiner, and a Prototype Branch for image-specific  
407 normal prototype construction. These designs are intended to alleviate three common issues  
408 in existing CLIP-based anomaly detection methods, namely semantic granularity mismatch  
409 across layers, inflexible fusion across feature levels, and interference from background or normal  
410 appearance variations.

411 Experimental results on multiple industrial defect and medical imaging datasets show that  
412 LaRP-CLIP achieves strong performance in both image-level detection and pixel-level localization.  
413 The qualitative results further show that the proposed framework produces cleaner anomaly

414 responses and more precise boundaries in challenging cases. These results suggest that LaRP-CLIP  
415 is a practical zero-shot solution for anomaly detection in scenarios where anomalous samples are  
416 limited or unavailable. Future work will consider extending the framework to multi-class anomaly  
417 settings and to more complex real-world environments.

#### 418 **Acknowledgments**

419 The authors gratefully acknowledge the anonymous reviewers and the editor for their  
420 thoughtful feedback and insightful suggestions, which substantially strengthened the clarity,  
421 precision, and rigor of this study.

#### 422 **Funding Statement**

423 This research was funded by the Key Research and Development Program of Zhejiang Province  
424 No. 2023C01141, the Science and Technology Innovation Community Project of Yangtze River Delta  
425 No. 23002410100.

426 This work was supported by the Open Research Fund of The State Key Laboratory of  
427 Blockchain and Data Security, Zhejiang University.

#### 428 **Author Contributions**

429 Xing Fang: Responsible for software implementation, experimental validation, and preparation  
430 of the original manuscript draft. Yuanfang Chen: Led the conceptualization and methodological  
431 design, provided project administration and supervision, contributed to manuscript revision  
432 and editing, and secured funding. Qiang Lin: Conducted data collection. Kun Yang: Offered  
433 methodological guidance, technical support, and constructive feedback during manuscript review.  
434 Gyu Myoung Lee: Provided conceptual guidance and supervision.

435 All authors have read and approved the final version of the manuscript.

#### 436 **Availability of Data and Materials**

437 All datasets used in this study are publicly accessible, and links to the original data sources are  
438 included in the manuscript for reference.

#### 439 **Ethics Approval**

440 Not applicable.

#### 441 **Conflicts of Interest**

442 The authors declare no known competing financial interests or personal relationships that  
443 could have influenced the work reported in this paper.

#### 444 **References**

445

- 446 1. Wu P, Pan C, Yan Y, Pang G, Yan Q, Wang P, et al. Deep learning for video anomaly detection: A  
447 review. IEEE Transactions on Neural Networks and Learning Systems. 2026. Available from: <https://ieeexplore.ieee.org/abstract/document/11322811/>.  
448

- 449 2. Huang H, Wang P, Pei J, Wang J, Alexanian S, Niyato D. Deep learning advancements in anomaly  
450 detection: A comprehensive survey. *IEEE Internet of Things Journal*. 2025. Available from: <https://ieeexplore.ieee.org/abstract/document/11071924/>.  
451
- 452 3. Li Z, Yan Y, Wang X, Ge Y, Meng L. A survey of deep learning for industrial visual anomaly detection.  
453 *Artificial Intelligence Review*. 2025;58(9):279. Available from: <https://link.springer.com/article/10.1007/s10462-025-11287-7>.  
454
- 455 4. Ammar MB, Mendoza A, Belkhir N, Manzanera A, Franchi G. Foundation models and Transformers  
456 for anomaly detection: A survey. *Information Fusion*. 2025:103517. Available from: <https://www.sciencedirect.com/science/article/pii/S1566253525005895>.  
457
- 458 5. Laghari AA, Khan AA, Ksibi A, Hajjej F, Kryvinska N, Almadhor A, et al. A novel and secure artificial  
459 intelligence enabled zero trust intrusion detection in industrial internet of things architecture. *Scientific*  
460 *Reports*. 2025;15(1):26843. Available from: <https://www.nature.com/articles/s41598-025-11738-9>.
- 461 6. Xie S, Wu X, Wang MY. Semi-patchcore: A novel two-staged method for semi-supervised anomaly  
462 detection and localization. *IEEE Transactions on Instrumentation and Measurement*. 2025;74:1-12.  
463 Available from: <https://ieeexplore.ieee.org/abstract/document/10835799>.
- 464 7. Li R, Ma H, Wang R, Song H, Zhou X, Wang L, et al. Application of unsupervised learning methods  
465 based on video data for real-time anomaly detection in wire arc additive manufacturing. *Journal of*  
466 *Manufacturing Processes*. 2025;143:37-55. Available from: [https://www.sciencedirect.com/science/](https://www.sciencedirect.com/science/article/pii/S1526612525003718)  
467 [article/pii/S1526612525003718](https://www.sciencedirect.com/science/article/pii/S1526612525003718).
- 468 8. Koren O, Koren M, Peretz O. A procedure for anomaly detection and analysis. *Engineering Applications*  
469 *of Artificial Intelligence*. 2023;117:105503. Available from: [https://www.sciencedirect.com/science/](https://www.sciencedirect.com/science/article/pii/S0952197622004936)  
470 [article/pii/S0952197622004936](https://www.sciencedirect.com/science/article/pii/S0952197622004936).
- 471 9. Nizam H, Zafar S, Lv Z, Wang F, Hu X. Real-time deep anomaly detection framework for multivariate  
472 time-series data in industrial IoT. *IEEE Sensors Journal*. 2022;22(23):22836-49. Available from: <https://ieeexplore.ieee.org/abstract/document/9915308>.  
473
- 474 10. Liu L, Li J, Lv J, Wang J, Zhao S, Lu Q. Privacy-preserving and secure industrial big data analytics: A  
475 survey and the research framework. *IEEE Internet of Things Journal*. 2024;11(11):18976-99. Available  
476 from: <https://ieeexplore.ieee.org/abstract/document/10399957>.
- 477 11. Aich A, Peng KC, Roy-Chowdhury AK. Cross-domain video anomaly detection without target domain  
478 adaptation. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*;  
479 2023. p. 2579-91. Available from: [https://openaccess.thecvf.com/content/WACV2023/html/Aich\\_](https://openaccess.thecvf.com/content/WACV2023/html/Aich_Cross-Domain_Video_Anomaly_Detection_Without_Target_Domain_Adaptation_WACV_2023_paper.html)  
480 [Cross-Domain\\_Video\\_Anomaly\\_Detection\\_Without\\_Target\\_Domain\\_Adaptation\\_WACV\\_2023\\_paper.](https://openaccess.thecvf.com/content/WACV2023/html/Aich_Cross-Domain_Video_Anomaly_Detection_Without_Target_Domain_Adaptation_WACV_2023_paper.html)  
481 [html](https://openaccess.thecvf.com/content/WACV2023/html/Aich_Cross-Domain_Video_Anomaly_Detection_Without_Target_Domain_Adaptation_WACV_2023_paper.html).
- 482 12. Hashemi MJ, Keller E, Tizpaz-Niari S. Detecting unseen anomalies in network systems by leveraging  
483 neural networks. *IEEE Transactions on Network and Service Management*. 2022;20(3):2515-28. Available  
484 from: <https://ieeexplore.ieee.org/abstract/document/9944024>.
- 485 13. Yang Z, Soltani I, Darve E. Anomaly detection with domain adaptation. In: *Proceedings of the*  
486 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*; 2023. p. 2958-67. Available  
487 from: [https://openaccess.thecvf.com/content/CVPR2023W/VAND/html/Yang\\_Anomaly\\_Detection\\_](https://openaccess.thecvf.com/content/CVPR2023W/VAND/html/Yang_Anomaly_Detection_With_Domain_Adaptation_CVPRW_2023_paper.html)  
488 [With\\_Domain\\_Adaptation\\_CVPRW\\_2023\\_paper.html](https://openaccess.thecvf.com/content/CVPR2023W/VAND/html/Yang_Anomaly_Detection_With_Domain_Adaptation_CVPRW_2023_paper.html).
- 489 14. Cho M, Kim T, Shim M, Wee D, Lee S. Towards multi-domain learning for generalizable video  
490 anomaly detection. *Advances in Neural Information Processing Systems*. 2024;37:50256-84. Available  
491 from: [https://openaccess.thecvf.com/content/CVPR2023W/VAND/html/Yang\\_Anomaly\\_Detection\\_](https://openaccess.thecvf.com/content/CVPR2023W/VAND/html/Yang_Anomaly_Detection_With_Domain_Adaptation_CVPRW_2023_paper.html)  
492 [With\\_Domain\\_Adaptation\\_CVPRW\\_2023\\_paper.html](https://openaccess.thecvf.com/content/CVPR2023W/VAND/html/Yang_Anomaly_Detection_With_Domain_Adaptation_CVPRW_2023_paper.html).
- 493 15. Mao W, Wang G, Kou L, Liang X. Deep domain-adversarial anomaly detection with one-class transfer  
494 learning. *IEEE/CAA Journal of Automatica Sinica*. 2023;10(2):524-46. Available from: [https://ieeexplore.](https://ieeexplore.ieee.org/abstract/document/10024159)  
495 [ieeexplore.](https://ieeexplore.ieee.org/abstract/document/10024159)  
496 [iee.org/abstract/document/10024159](https://ieeexplore.ieee.org/abstract/document/10024159).
- 497 16. Li A, Qiu C, Kloft M, Smyth P, Rudolph M, Mandt S. Zero-shot anomaly detection via batch  
498 normalization. *Advances in Neural Information Processing Systems*. 2023;36:40963-93. Available  
499 from: [https://proceedings.neurips.cc/paper\\_files/paper/2023/hash/8078e8c3055303a884ffae2d3ea0](https://proceedings.neurips.cc/paper_files/paper/2023/hash/8078e8c3055303a884ffae2d3ea00338-Abstract-Conference.html)  
0338-Abstract-Conference.html.

- 500 17. Peng Y, Lin X, Ma N, Du J, Liu C, Liu C, et al. Sam-lad: Segment anything model meets zero-shot  
501 logic anomaly detection. *Knowledge-Based Systems*. 2025;314:113176. Available from: [https://www.  
502 sciencedirect.com/science/article/pii/S0950705125002230](https://www.sciencedirect.com/science/article/pii/S0950705125002230).
- 503 18. Radford A, Kim JW, Hallacy C, Ramesh A, Goh G, Agarwal S, et al. Learning transferable visual models  
504 from natural language supervision. In: *International conference on machine learning*. PmlR; 2021. p.  
505 8748-63. Available from: <https://proceedings.mlr.press/v139/radford21a>.
- 506 19. Jeong J, Zou Y, Kim T, Zhang D, Ravichandran A, Dabeer O. Winclip: Zero-/few-shot anomaly  
507 classification and segmentation. In: *Proceedings of the IEEE/CVF conference on computer vision and  
508 pattern recognition*; 2023. p. 19606-16. Available from: [https://openaccess.thecvf.com/content/CVPR20  
509 023/html/Jeong\\_WinCLIP\\_Zero-Few-Shot\\_Anomaly\\_Classification\\_and\\_Segmentation\\_CVPR\\_2023  
510 \\_paper.html](https://openaccess.thecvf.com/content/CVPR2023/html/Jeong_WinCLIP_Zero-Few-Shot_Anomaly_Classification_and_Segmentation_CVPR_2023_paper.html).
- 511 20. Zhou Q, Pang G, Tian Y, He S, Chen J. Anomalyclip: Object-agnostic prompt learning for zero-shot  
512 anomaly detection. In: *International conference on learning representation*; 2023. Available from:  
513 <https://arxiv.org/abs/2310.18961>.
- 514 21. Kim D, Park C, Cho S, Lim H, Kang M, Lee J, et al. Generalizing CLIP Prompts for Zero-shot Anomaly  
515 Detection. *Pattern Recognition*. 2026:113406. Available from: [https://www.sciencedirect.com/science/  
516 article/abs/pii/S0031320326003717](https://www.sciencedirect.com/science/article/abs/pii/S0031320326003717).
- 517 22. Chen X, Han Y, Zhang J. A zero-/few-shot anomaly classification and segmentation method for CVPR  
518 2023 (VAND) workshop challenge tracks 1 & 2. 1st Place on Zero-shot AD and 4th Place on Few-shot AD.  
519 2023;2305(17382):2. Available from: <https://arxiv.org/abs/2305.17382>.
- 520 23. Zou Y, Jeong J, Pemula L, Zhang D, Dabeer O. Spot-the-difference self-supervised pre-training for  
521 anomaly detection and segmentation. In: *European conference on computer vision*. Springer; 2022. p.  
522 392-408. Available from: [https://link.springer.com/chapter/10.1007/978-3-031-20056-4\\_23](https://link.springer.com/chapter/10.1007/978-3-031-20056-4_23).
- 523 24. Ma W, Zhang X, Yao Q, Tang F, Wu C, Li Y, et al. Aa-clip: Enhancing zero-shot anomaly detection via  
524 anomaly-aware clip. In: *Proceedings of the Computer Vision and Pattern Recognition Conference*; 2025.  
525 p. 4744-54. Available from: [https://openaccess.thecvf.com/content/CVPR2025/html/Ma\\_AA-CLIP\\_  
526 Enhancing\\_Zero-Shot\\_Anomaly\\_Detection\\_via\\_Anomaly-Aware\\_CLIP\\_CVPR\\_2025\\_paper.html](https://openaccess.thecvf.com/content/CVPR2025/html/Ma_AA-CLIP_Enhancing_Zero-Shot_Anomaly_Detection_via_Anomaly-Aware_CLIP_CVPR_2025_paper.html).
- 527 25. Ruff L, Kauffmann JR, Vandermeulen RA, Montavon G, Samek W, Kloft M, et al. A unifying review  
528 of deep and shallow anomaly detection. *Proceedings of the IEEE*. 2021;109(5):756-95. Available from:  
529 <https://ieeexplore.ieee.org/abstract/document/9347460/>.
- 530 26. Gong D, Liu L, Le V, Saha B, Mansour MR, Venkatesh S, et al. Memorizing normality to detect  
531 anomaly: Memory-augmented deep autoencoder for unsupervised anomaly detection. In: *Proceedings  
532 of the IEEE/CVF international conference on computer vision*; 2019. p. 1705-14. Available from:  
533 [https://openaccess.thecvf.com/content\\_ICCV\\_2019/html/Gong\\_Memorizing\\_Normality\\_to\\_Detect\\_  
534 Anomaly\\_Memory-Augmented\\_Deep\\_Autoencoder\\_for\\_Unsupervised\\_ICCV\\_2019\\_paper.html](https://openaccess.thecvf.com/content_ICCV_2019/html/Gong_Memorizing_Normality_to_Detect_Anomaly_Memory-Augmented_Deep_Autoencoder_for_Unsupervised_ICCV_2019_paper.html).
- 535 27. Schlegl T, Seeböck P, Waldstein SM, Langs G, Schmidt-Erfurth U. f-AnoGAN: Fast unsupervised anomaly  
536 detection with generative adversarial networks. *Medical image analysis*. 2019;54:30-44. Available from:  
537 <https://www.sciencedirect.com/science/article/pii/S1361841518302640>.
- 538 28. Ruff L, Vandermeulen R, Goernitz N, Deecke L, Siddiqui SA, Binder A, et al. Deep one-class classification.  
539 In: *International conference on machine learning*. PMLR; 2018. p. 4393-402. Available from: <https://proceedings.mlr.press/v80/ruff18a.html>.
- 541 29. Ren J, Tang T, Jia H, Xu Z, Fayek H, Li X, et al. Foundation models for anomaly detection: Vision and  
542 challenges. *AI Magazine*. 2025;46(4):e70045. Available from: [https://onlinelibrary.wiley.com/doi/abs/  
543 10.1002/aaai.70045](https://onlinelibrary.wiley.com/doi/abs/10.1002/aaai.70045).
- 544 30. Li Y, Goodge A, Liu F, Foo CS. Promptad: Zero-shot anomaly detection using text prompts. In:  
545 *Proceedings of the IEEE/CVF winter conference on applications of computer vision*; 2024. p. 1093-102.  
546 Available from: [https://openaccess.thecvf.com/content/WACV2024/html/Li\\_PromptAD\\_Zero-Shot\\_  
547 Anomaly\\_Detection\\_Using\\_Text\\_Prompts\\_WACV\\_2024\\_paper.html](https://openaccess.thecvf.com/content/WACV2024/html/Li_PromptAD_Zero-Shot_Anomaly_Detection_Using_Text_Prompts_WACV_2024_paper.html).
- 548 31. Chen Z, Zhao Z, Guo J, Li J, Huang Z. SVIP: Semantically Contextualized Visual Patches for Zero-Shot  
549 Learning. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*; 2025.

- 550 Available from: [https://openaccess.thecvf.com/content/ICCV2025/papers/Chen\\_SVIP\\_Semantically\\_](https://openaccess.thecvf.com/content/ICCV2025/papers/Chen_SVIP_Semantically_)  
551 [Contextualized\\_Visual\\_Patches\\_for\\_Zero-Shot\\_Learning\\_ICCV\\_2025\\_paper.pdf](https://openaccess.thecvf.com/content/ICCV2025/papers/Chen_SVIP_Semantically_).
- 552 32. Bergmann P, Fauser M, Sattlegger D, Steger C. MVTEC AD—A comprehensive real-world dataset for  
553 unsupervised anomaly detection. In: Proceedings of the IEEE/CVF conference on computer vision and  
554 pattern recognition; 2019. p. 9592-600. Available from: [https://openaccess.thecvf.com/content\\_CVPR\\_](https://openaccess.thecvf.com/content_CVPR_2019/html/Bergmann_MVTEC_AD_--_A_Comprehensive_Real-World_Dataset_for_Unsupervised_)  
555 [2019/html/Bergmann\\_MVTEC\\_AD\\_--\\_A\\_Comprehensive\\_Real-World\\_Dataset\\_for\\_Unsupervised\\_](https://openaccess.thecvf.com/content_CVPR_2019/html/Bergmann_MVTEC_AD_--_A_Comprehensive_Real-World_Dataset_for_Unsupervised_)  
556 [Anomaly\\_CVPR\\_2019\\_paper.html](https://openaccess.thecvf.com/content_CVPR_2019/html/Bergmann_MVTEC_AD_--_A_Comprehensive_Real-World_Dataset_for_Unsupervised_).
- 557 33. Jezek S, Jonak M, Burget R, Dvorak P, Skotak M. Deep learning-based defect detection of metal parts:  
558 evaluating current methods in complex conditions. In: 2021 13th International congress on ultra modern  
559 telecommunications and control systems and workshops (ICUMT). IEEE; 2021. p. 66-71. Available from:  
560 <https://ieeexplore.ieee.org/abstract/document/9631567/>.
- 561 34. Mishra P, Verk R, Fornasier D, Piciarelli C, Foresti GL. VT-ADL: A vision transformer network for  
562 image anomaly detection and localization. arXiv preprint arXiv:2104.10036. 2021. Available from:  
563 <https://arxiv.org/abs/2104.10036>.
- 564 35. Salehi M, Sadjadi N, Baselizadeh S, Rohban MH, Rabiee HR. Multiresolution knowledge distillation  
565 for anomaly detection. In: Proceedings of the IEEE/CVF conference on computer vision and pattern  
566 recognition; 2021. p. 14902-12. Available from: [https://openaccess.thecvf.com/content/CVPR2021/](https://openaccess.thecvf.com/content/CVPR2021/html/Salehi_Multiresolution_Knowledge_Distillation_for_Anomaly_Detection_CVPR_2021_paper.html)  
567 [html/Salehi\\_Multiresolution\\_Knowledge\\_Distillation\\_for\\_Anomaly\\_Detection\\_CVPR\\_2021\\_paper.](https://openaccess.thecvf.com/content/CVPR2021/html/Salehi_Multiresolution_Knowledge_Distillation_for_Anomaly_Detection_CVPR_2021_paper.html)  
568 [html](https://openaccess.thecvf.com/content/CVPR2021/html/Salehi_Multiresolution_Knowledge_Distillation_for_Anomaly_Detection_CVPR_2021_paper.html).
- 569 36. Hamada A. Br35H :: Brain Tumor Detection 2020. IEEE Dataport; 2025. Available from: [https://dx.doi.](https://dx.doi.org/10.21227/t5k6-5r54)  
570 [org/10.21227/t5k6-5r54](https://dx.doi.org/10.21227/t5k6-5r54).
- 571 37. Codella NC, Gutman D, Celebi ME, Helba B, Marchetti MA, Dusza SW, et al. Skin lesion analysis toward  
572 melanoma detection: A challenge at the 2017 international symposium on biomedical imaging (isbi),  
573 hosted by the international skin imaging collaboration (isic). In: 2018 IEEE 15th international symposium  
574 on biomedical imaging (ISBI 2018). IEEE; 2018. p. 168-72. Available from: [https://ieeexplore.ieee.org/](https://ieeexplore.ieee.org/abstract/document/8363547/)  
575 [abstract/document/8363547/](https://ieeexplore.ieee.org/abstract/document/8363547/).
- 576 38. Tajbakhsh N, Gurudu SR, Liang J. Automated polyp detection in colonoscopy videos using shape  
577 and context information. IEEE transactions on medical imaging. 2015;35(2):630-44. Available from:  
578 <https://ieeexplore.ieee.org/abstract/document/7294676/>.
- 579 39. Gong H, Chen G, Wang R, Xie X, Mao M, Yu Y, et al. Multi-task learning for thyroid nodule segmentation  
580 with thyroid region prior. In: 2021 IEEE 18th international symposium on biomedical imaging (ISBI).  
581 IEEE; 2021. p. 257-61. Available from: <https://ieeexplore.ieee.org/abstract/document/9434087/>.
- 582 40. Zhou K, Yang J, Loy CC, Liu Z. Learning to prompt for vision-language models. International journal of  
583 computer vision. 2022;130(9):2337-48. Available from: [https://link.springer.com/article/10.1007/s11263](https://link.springer.com/article/10.1007/s11263-022-01653-1)  
584 [-022-01653-1](https://link.springer.com/article/10.1007/s11263-022-01653-1).