

Toward Generalist Anomaly Detection via In-context Residual Learning with Few-shot Sample Prompts

Jiawen Zhu and Guansong Pang *

School of Computing and Information Systems, Singapore Management University

Abstract

This paper explores the problem of Generalist Anomaly Detection (GAD), aiming to train one single detection model that can generalize to detect anomalies in diverse datasets from different application domains without any further training on the target data. Some recent studies have shown that large pre-trained Visual-Language Models (VLMs) like CLIP have strong generalization capabilities on detecting industrial defects from various datasets, but their methods rely heavily on handcrafted text prompts about defects, making them difficult to generalize to anomalies in other applications, e.g., medical image anomalies or semantic anomalies in natural images. In this work, we propose to **train a GAD model with few-shot normal images as sample prompts for AD on diverse datasets on the fly**. To this end, we introduce a novel approach that learns an in-context residual learning model for GAD, termed **InCTRL**. It is trained on an auxiliary dataset to discriminate anomalies from normal samples based on a holistic evaluation of the residuals between query images and few-shot normal sample prompts. Regardless of the datasets, per definition of anomaly, larger residuals are expected for anomalies than normal samples, thereby enabling InCTRL to generalize across different domains without further training. Comprehensive experiments on nine AD datasets are performed to establish a GAD benchmark that encapsulate the detection of industrial defect anomalies, medical anomalies, and semantic anomalies in both one-vs-all and multi-class setting, on which InCTRL is the best performer and significantly outperforms state-of-the-art competing methods. Code is available at <https://github.com/mala-lab/InCTRL>.

1. Introduction

Anomaly Detection (AD) is a crucial computer vision task that aims to detect samples that substantially deviate from

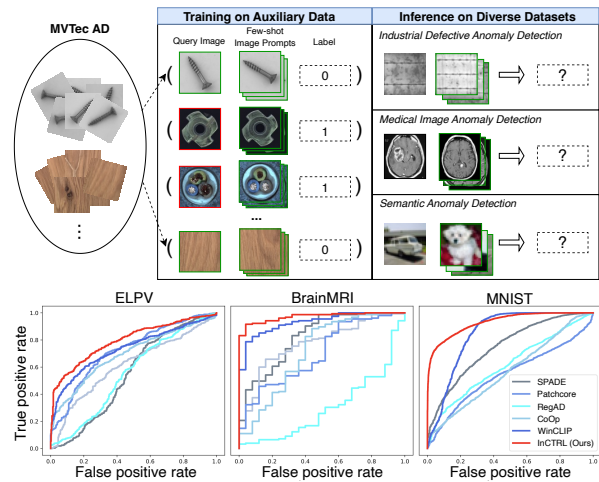


Figure 1. **Top**: An illustration of InCTRL: a one-for-all model using few-shot normal images as sample prompts. **Bottom**: AUROC curves of InCTRL and competing few-shot methods on three different application datasets without any training on the target data.

the majority of samples in a dataset, due to its broad real-life applications such as industrial inspection, medical imaging analysis, and scientific discovery, etc. [12, 37]. **Current AD paradigms are focused on individually building one model on the training data**, e.g., a set of anomaly-free samples, of each target dataset, such as data reconstruction approach [2, 22, 32, 39, 44, 61, 64–66, 68–70], one-class classification [6, 15, 42, 49, 67], and knowledge distillation approach [8, 11, 19, 43, 52, 53, 71].

Although these approaches have shown remarkable detection performance on various AD benchmarks, **they require the availability of large training data and the skillful training of the detection model per dataset**. Thus, they become infeasible in application scenarios where **training on the target dataset is not allowed** due to either data privacy issues, e.g., arising from using those data in training the models due to machine unlearning [63], or unavailability of large-scale training data in the deployment of new applications. To tackle these challenges, this paper explores

*Corresponding author: G. Pang (gspang@smu.edu.sg)

the problem of learning **Generalist Anomaly Detection** (GAD) models, *aiming to train one single detection model that can generalize to detect anomalies in diverse datasets from different application domains without any training on the target data.*

Being pre-trained on web-scale image-text data, large Visual-Language Models (VLMs) like CLIP [40] have exhibited superior generalization capabilities in recent years, achieving accurate visual recognition across different datasets without any fine-tuning or adaptation on the target data. More importantly, some very recent studies (e.g., WinCLIP [25]) show that these VLMs can also be utilized to achieve remarkable generalization on different defect detection datasets. Nevertheless, a significant limitation of these models is their dependency on a large set of manually crafted prompts specific to defects. This reliance restricts their applicability, making it challenging to extend their use to detecting anomalies in other data domains, e.g., medical image anomalies [10, 20, 43, 50, 51] or semantic anomalies in one-vs-all or multi-class settings [11, 42].

To address this problem, we propose to train a GAD model that aims to utilize few-shot normal images from any target dataset as sample prompts for supporting GAD on the fly, as illustrated in Figure 1(Top). The few-shot setting is motivated by the fact that it is often easy to obtain few-shot normal images in real-world applications. Furthermore, these few-shot samples are not used for model training/tuning; they are just used as sample prompts for enabling the anomaly scoring of test images during inference. This formulation is fundamentally different from current few-shot AD methods [5, 23, 45, 46, 57, 58, 62] that use these target samples and their extensive augmented versions to train the detection model, which can lead to an overfitting of the target dataset and fail to generalize to other datasets, as shown in Figure 1(Bottom).

We then introduce a GAD approach, the first of its kind, that learns an in-context residual learning model based on CLIP, termed InCTRL. It trains a GAD model to discriminate anomalies from normal samples by learning to identify the residuals/discrepancies between query images and a set of few-shot normal images from auxiliary data. The few-shot normal images, namely in-context sample prompts, serve as prototypes of normal patterns. When comparing with the features of these normal patterns, per definition of anomaly, a larger residual is typically expected for anomalies than normal samples in datasets of different domains, so the learned in-context residual model can generalize to detect diverse types of anomalies across the domains.

To capture the residuals better, InCTRL models the in-context residuals at both the image and patch levels, gaining an in-depth in-context understanding of what constitutes an anomaly. Further, our in-context residual learning can also enable a seamless incorporation of normal/abnormal text

prompt-guided prior knowledge into the detection model, providing an additional strength for the detection from the text-image-aligned semantic space.

Accordingly, we make the following main contributions.

- We introduce a GAD task to evaluate the generalization capability of AD methods in identifying anomalies across various scenarios without needing to training/tuning on the target datasets. To the best of our knowledge, this is the first study dedicated to a generalist approach to anomaly detection, encompassing industrial defects, medical anomalies, and semantic anomalies.
- We then propose an in-context residual learning framework for GAD, called InCTRL. It is designed to distinguish anomalies from normal samples by detecting residuals between test images and in-context few-shot normal sample prompts from any target dataset on the fly. InCTRL is optimized on auxiliary data to achieve the one-model-for-all goal, i.e., one model for AD on diverse datasets without any training on target data.
- Comprehensive experiments on nine diverse AD datasets are performed to establish a GAD benchmark that encapsulates three types of popular AD tasks, including industrial defect anomaly detection, medical image anomaly detection, and semantic anomaly detection under both one-vs-all and multi-class settings. Our results show that InCTRL significantly outperforms state-of-the-art competing methods.

2. Related Work

2.1. Anomaly Detection

Anomaly Detection. Existing AD approaches typically rely on unsupervised learning due to the scarcity of anomaly data. Numerous methods have been introduced. One-class classification methods [6, 15, 42, 49, 67] focus on compactly describing normal data with support vectors. Reconstruction-based methods [2, 22, 32, 39, 44, 61, 64–66, 68–70] train models to reconstruct normal images, where anomalies are identified by higher reconstruction errors. Distance-based methods [16, 17, 41] determine anomalies based on the distance between test image embeddings and normal reference embeddings from stored training data. Knowledge distillation methods [8, 11, 19, 43, 52, 53, 71] focus on distilling normal patterns from pre-trained models and detect anomalies based on the difference between distilled and original features. The above approaches are designed to fit on target dataset for AD, i.e., one model for one dataset. We aim for a one-model-for-all setting. A relevant research line is to tackle the AD problem under domain or distribution shift [1, 11, 20, 34, 66, 74], but they generally assume a large domain relevance on the source and target data. Additionally, there have been a number of concurrent studies leveraging VLMs for AD [59, 60, 73], but they ad-

dress a different setting from ours, *e.g.*, weakly-supervised AD [59, 60] or zero-shot AD [73].

Few-shot Anomaly Detection (FSAD), FSAD is designed to identify anomalies using only a limited number of normal samples from target datasets. Traditional FSAD research focuses on modeling the normal distribution of these few normal samples to detect anomalies [5, 23, 30, 45, 46, 57, 58, 62]. However, these methods often cannot generalize to new domains, as they generally require re-training or fine-tuning with normal data from the target datasets.

Distance-based approaches such as SPADE [16], PaDiM [17] and PatchCore [41] present a solution to address this problem by making full use of available pre-trained representations of the few-shot samples to calculate distance-based anomaly scores without training. Recently, **RegAD** [23] is designed as a model that operates without the need for training or fine-tuning on new data for the FSAD task, but it requires domain relevance between training and test data to work well. WinCLIP [25] pioneers the application of large Visual-Language Models (VLM) on zero-shot and few-shot anomaly detection tasks by processing images through multi-scale window movements and text prompting to CLIP. Without adapting CLIP to the AD task, WinCLIP gains impressive zero-shot detection performance on defect datasets using its handcrafted text prompts, but it fails to work well when the text prompts cannot capture the required anomaly semantics, making it difficult to generalize well to diverse anomaly detection tasks.

2.2. In-Context Learning

In-context learning is an innovative approach that helps enhance the performance of Large Language Models (LLMs) in Natural Language Processing (NLP) [3, 9, 21], which leverages minimal in-context prompts to adapt LLMs to novel tasks effectively.

Recently, several studies [13, 14, 26, 33, 54] attempt to apply in-context learning to vision tasks by converting vision problems to NLP ones using the language or specially-designed discrete tokens as the task prompts. On the other hand, Amir *et al.* [4] introduce a novel approach for in-context visual prompting by treating a spectrum of vision tasks as grid in-painting problems. Similarly, Painter [55, 56] then proposes to perform masked image in-painting. However, these methods focus more on task-level generalization, so they are not applicable to the AD task which focuses more on the instance-level discrepancy.

Our work redesign in-context learning for GAD. We redefine image prompts as dataset-specific normal patterns, rather than as an instruction for particular tasks. By capturing the in-context residual between the query image and few-shot normal prompts, our model can gain a cohesive understanding of diverse anomalies, enabling remarkable generalized detection performance for GAD.

3. InCTRL: In-Context Residual Learning

3.1. Problem Statement

The objective of GAD is to train a single AD model that works well for detecting anomalies on test datasets from diverse application domains without any training on the target data. Thus, the training set is assumed to be drawn from different distributions from the test sets. Formally, let $\mathcal{D}_{train} = \{X_{train}, Y_{train}\}$ be an *auxiliary* training dataset with normal and anomaly class labels, where $X_{train} = \{x_i\}_{i=1}^N$ consists of N normal and anomalous images and $Y_{train} = \{y_i\}_{i=1}^N$, with $y_i = 0$ indicates normal and $y_i = 1$ signifies abnormal. A collection of test sets, $\mathcal{T} = \{\mathcal{D}_{test}^1, \mathcal{D}_{test}^2, \dots, \mathcal{D}_{test}^M\}$ with $\mathcal{D}_{test}^j = \{X_{test}^j, Y_{test}^j\}$, from M different application domains with various types of anomalies is given. The test sets are drawn from a distribution different from that of \mathcal{D}_{train} . Then the goal is to train a generalist anomaly scoring function: $\mathcal{D}_{train} \rightarrow \mathbb{R}$ so that it assigns larger anomaly scores to the anomalous samples than to the normal ones from any test dataset in \mathcal{T} . In the context of GAD with few-shot normal samples, a small set of a few normal images randomly drawn from the target domain, $\mathcal{P} = \{p_1, p_2, \dots, p_K\}$ where K is typically a small number, *e.g.*, $K \ll N$, is available during inference, but \mathcal{P} is not available in any way during the training of the generalist detection model

3.2. Overview of Our Approach InCTRL

Our approach InCTRL is designed to effectively model the in-context residual between a query image and a set of few-shot normal images as sample prompts, utilizing the generalization capabilities of CLIP to detect unusual residuals for anomalies from different application domains. CLIP is a VLM consisting of a text encoder $f_t(\cdot)$ and a visual encoder $f_v(\cdot)$, with the image and text representations from these encoders well aligned by pre-training on web-scale text-image data. InCTRL is optimized using auxiliary data \mathcal{D}_{train} via an in-context residual learning in the image encoder, with the learning augmented by text prompt-guided prior knowledge from the text encoder.

To be more specific, as illustrated in Fig. 2, we first simulate an in-context learning example that contains one query image x and a set of few-shot normal sample prompts \mathcal{P}' , both of which are randomly sampled from the auxiliary data \mathcal{D}_{train} . Through the visual encoder, we then perform multi-layer patch-level and image-level residual learning to respectively capture local and global discrepancies between the query and few-shot normal sample prompts (Secs. 3.3 and 3.4). Further, our model allows a seamless incorporation of normal and abnormal text prompts-guided prior knowledge from the text encoder based on the similarity between these textual prompt embeddings and the query images (Sec. 3.5). The training of InCTRL is to optimize a

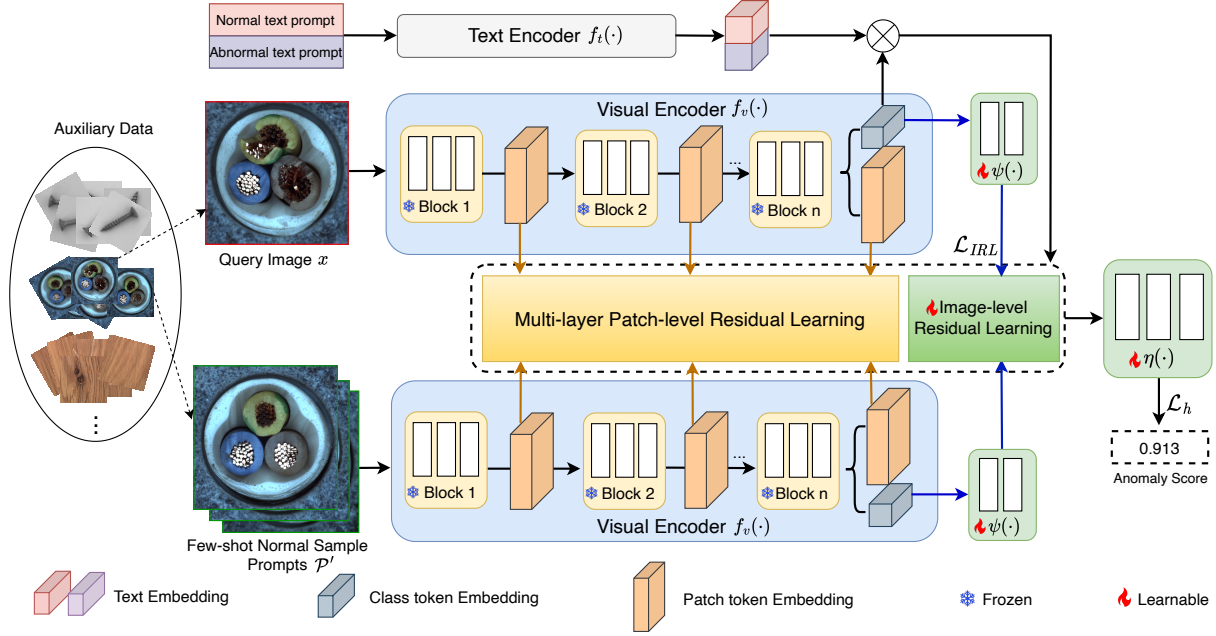


Figure 2. Overview of the training of InCTRL. It firstly simulates in-context learning scenarios using a query image and a few-shot normal sample prompts randomly drawn from the auxiliary training data. Then it performs multi-layer patch-level and image-level residual learning to capture both local and global residuals between the query image and the normal prompts. Lastly, those residual information, combined with text prompts-guided prior knowledge from the text encoder, is utilized for a holistic anomaly score learning.

few projection/adaptation layers attached to the visual encoder to learn a larger anomaly score for anomaly samples than normal samples in \mathcal{D}_{train} , with the original parameters in both encoders frozen; during inference, a test image, together with the few-shot normal image prompts from the target dataset and the text prompts, is put forward through our adapted CLIP-based GAD network, whose output is the anomaly score for the test image (Sec. 3.6). Below we present these modules in detail.

3.3. Multi-Layer Patch-Level Residual Learning

To effectively capture fine-grained in-context residuals between the query image and the normal image prompts, we introduce a multi-layer patch-level residual learning component in InCTRL. Typically, the CLIP visual encoder comprises a series of block layers. From the bottom to the top of layers, the visual encoder gradually learns the visual patterns at different levels of abstraction [40]. Thus, this component is designed to model patch-level in-context residuals from the patch token embeddings obtained from the multiple levels of the blocks within the visual encoder.

To be specific, assuming the visual encoder consists of n blocks, for a given set of few-shot normal sample prompts \mathcal{P}' and a training query image x , we extract a series of patch token embedding maps $\{T_x^l\}_{l=1}^n$ and $\{T_{x'}^l\}_{l=1}^n$ where $T_{(\cdot)}^l \in \mathbb{R}^{h \times w \times d}$ and $x' \in \mathcal{P}'$, with h , w , and d be the height, width, and dimension of the feature map T respectively. At each layer l , the patch-level in-context residuals

are captured by distances between the embeddings of the query token and the image prompt token across all image prompts in \mathcal{P}' . Formally, for the query image x , its multi-layer patch-level in-context residuals at layer l are modeled by a residual map $\mathbf{M}_x^l \in \mathbb{R}^{h \times w}$, where the residual value of each patch of x is calculated based on its patch embedding and the nearest patch embedding of all images in \mathcal{P}' as:

$$\mathbf{M}_x^l(i, j) = 1 - \langle T_x^l(i, j), h(T_x^l(i, j) | \mathcal{P}') \rangle, \quad (1)$$

where $h(T_x^l(i, j) | \mathcal{P}')$ returns the embedding of the patch token that is most similar to $T_x^l(i, j)$ among all image patches in \mathcal{P}' , and $\langle \cdot \rangle$ is the cosine similarity function. The final patch-level residual map $\mathbf{M}_x \in \mathbb{R}^{h \times w}$ is averaged over n layer-wise residual maps:

$$\mathbf{M}_x = \frac{1}{n} \sum_{l=1}^n \mathbf{M}_x^l. \quad (2)$$

Each residual value in \mathbf{M}_x is similar to a nearest-neighbor-distance anomaly score of the query patch to the image patch set in \mathcal{P}' . As shown in prior studies [16, 17, 35, 36, 41], such distance-based anomaly scores can effectively discriminate anomalies from normal samples. Thus, the resulting residual map \mathbf{M}_x provides a feature set of collective anomaly-discriminative power at multi-layer resolutions for the subsequent anomaly score learning in InCTRL.

3.4. Image-level Residual Learning

In addition to the discriminative power at the local patch-level residuals, the global discriminative information at the image level is also significant and serve as complementary knowledge to the patch-level features.

Hence, we introduce an image-level residual learning component to capture the higher-level discrepancies between x and \mathcal{P}' . Intuitively, the class token embedding from the last block of the visual encoder is used as the feature input, as it captures the most image-level discriminative information due to the bottom-up abstraction of information in the visual encoder. However, it is important to note that CLIP was originally designed for classification tasks, focusing on the semantic of the objects in the scenery, which does not align well with the anomaly detection task in which both normal and abnormal samples are often from the same class of object. To reconcile this, we include an adapter layer $\psi(\cdot)$, parameterized by Θ_ψ , to adapt the image representations further to anomaly detection, and thus, we learn the image-level residuals based on the adapted image features. Further, the prototypical features of the few-shot sample prompts, rather than the features of individual sample, are used to learn the in-context residuals, since they help capture more representative features of normal patterns.

Specifically, let $f_v(x) \in \mathbb{R}^{d'}$ be the class token embedding of input x in the visual encoder, we first compute the prototype of the feature maps of the image prompts in \mathcal{P}' :

$$\mathbf{I}_p = \frac{1}{K} \sum_{x'_k \in \mathcal{P}'} \psi(f_v(x'_k); \Theta_\psi), \quad (3)$$

where $\mathbf{I}_p \in \mathbb{R}^{d'}$. Then let $\mathbf{I}_x = \psi(f_v(x); \Theta_\psi)$ be the adapted features of the query image x , the in-context image-level residual features \mathbf{F}_x for x are obtained by performing element-wise subtraction between two feature maps:

$$\mathbf{F}_x = \mathbf{I}_x \ominus \mathbf{I}_p, \quad (4)$$

where \ominus denotes element-wise subtraction. Subsequently, these in-context residual features are fed to an image-level anomaly classification learner $\eta : \mathbf{F}_x \rightarrow \mathbb{R}$, parameterized by Θ_η which is optimized by the binary classification loss:

$$\mathcal{L}_{IRL} = \frac{1}{N} \sum_{x \in X_{train}} \mathcal{L}_b(\eta(\mathbf{F}_x; \Theta_\eta), y_x), \quad (5)$$

where \mathcal{L}_b is a binary classification loss. Focal loss [31] is used by default in our model.

3.5. Fusing Text Prompt-based Prior Knowledge

The above two components are focused on residual learning based on the visual encoder. InCTRL also allows easy incorporation of text-prompt-guided prior knowledge about normality and abnormality from the text encoder of

CLIP. This helps InCTRL leverage the normal and abnormal semantics hidden in the CLIP’s pre-trained image-text-aligned embedding space for GAD. Motivated by this, InCTRL exploits the text encoder to extract text prompt-guided discriminative features. Since the text prompts designed in WinCLIP [25] show remarkable detection performance, InCTRL adopts the same text prompt templates and its ensemble strategy, including both state and template-level text prompts. At the state level, generic text descriptions are employed to differentiate between normal and abnormal objects, whereas the template level provides a list of specific prompts tailored for anomaly detection (see Appendix B.3 for detailed text prompts used).

It should be noted that, unlike WinCLIP that uses these text prompts to directly compute the anomaly score, InCTRL utilizes them to extract text-prompt-guided features for complementing the patch- and image-level residual features obtained through the visual encoder.

Specifically, let \mathcal{P}_t^n be the set of text prompts for the normal class, we use the prototype of the text prompt embeddings to provide a representative embedding of the normal text prompts $\mathbf{F}_n = \frac{1}{|\mathcal{P}_t^n|} \sum_{p_i \in \mathcal{P}_t^n} f_t(p_i)$ where $p_i \in \mathbb{R}^{d'}$; similarly we can obtain the prototype embedding for the abnormality text prompt set \mathcal{P}_t^a by $\mathbf{F}_a = \frac{1}{|\mathcal{P}_t^a|} \sum_{p_j \in \mathcal{P}_t^a} f_t(p_j)$. Then, InCTRL extracts an AD-oriented discriminative feature based on the similarity between the query image x and the two prototypes of the text prompts:

$$s_a(x) = \frac{\exp(\mathbf{F}_a^\top f_v(x))}{\exp(\mathbf{F}_n^\top f_v(x)) + \exp(\mathbf{F}_a^\top f_v(x))}, \quad (6)$$

where $[\cdot]^\top$ denotes a transpose operation, and $s_a(x)$ is the probability of the input x being classified as abnormal.

3.6. Training and Inference

In-Context Residual Learning. During training, InCTRL performs a holistic residual learning that synthesizes both patch-level and image-level residual information, augmented by the text prompt-guided features. The holistic in-context residual map of a query image x is defined as:

$$\mathbf{M}_x^+ = \mathbf{M}_x \oplus s_i(x) \oplus s_a(x), \quad (7)$$

where $s_i(x) = \eta(\mathbf{F}_x; \Theta_\eta)$ is an anomaly score based on the image-level residual map \mathbf{F}_x and \oplus denotes an element-wise addition. InCTRL then devises a holistic anomaly scoring function ϕ , parameterized by Θ_ϕ , based on \mathbf{M}_x^+ , and defines the final anomaly score as:

$$s(x) = \phi(\mathbf{M}_x^+; \Theta_\phi) + \alpha s_p(x), \quad (8)$$

where $\phi(\mathbf{M}_x^+; \Theta_\phi)$ performs a holistic anomaly scoring using patch-, image-level and text prompt-guided features, while $s_p(x) = \max(\mathbf{M}_x)$ is a maximum residual score-based fine-grained anomaly score at the image patch level.

Setup	Methods	Industrial Defects					Medical Anomalies		Semantic Anomalies			
		ELPV	SDD	AITEX	VisA	MVTec AD	BrainMRI	HeadCT	One-vs-all		Multi-class	
									MNIST	CIFAR-10	MNIST	CIFAR-10
2-shot	Baseline (0-shot)	0.733±0.000	0.946±0.000	0.733±0.000	0.781±0.000	0.912±0.000	0.926±0.000	0.900±0.000	0.678±0.000	0.924±0.000	0.620±0.000	0.900±0.000
	SPADE	0.517±0.012	0.729±0.041	0.727±0.004	0.795±0.045	0.817±0.054	0.754±0.048	0.645±0.034	0.779±0.024	0.823±0.014	0.595±0.060	0.655±0.042
	PaDiM	0.594±0.083	0.721±0.015	0.784±0.028	0.680±0.042	0.785±0.025	0.657±0.122	0.595±0.036	-	-	-	-
	Patchcore	0.716±0.031	0.902±0.006	0.739±0.017	0.817±0.028	0.858±0.034	0.706±0.009	0.736±0.096	0.756±0.004	0.602±0.009	0.603±0.009	0.703±0.008
	RegAD	0.571±0.016	0.499±0.008	0.564±0.072	0.557±0.053	0.640±0.047	0.449±0.129	0.602±0.018	0.525±0.030	0.534±0.005	0.608±0.026	0.695±0.002
	CoOp	0.762±0.011	0.897±0.006	0.687±0.062	0.806±0.023	0.858±0.016	0.725±0.020	0.811±0.003	0.557±0.006	0.527±0.011	0.612±0.007	0.393±0.009
	WinCLIP	0.726±0.020	0.942±0.006	0.726±0.055	0.842±0.024	0.931±0.019	0.934±0.012	0.915±0.015	0.810±0.008	0.925±0.001	0.632±0.000	0.914±0.005
	Ours (InCTRL)	0.839±0.003	0.972±0.011	0.761±0.029	0.858±0.022	0.940±0.015	0.973±0.027	0.929±0.025	0.892±0.009	0.935±0.002	0.635±0.010	0.924±0.005
4-shot	SPADE	0.537±0.013	0.731±0.020	0.718±0.011	0.811±0.040	0.828±0.044	0.759±0.070	0.624±0.012	0.810±0.009	0.836±0.006	0.588±0.041	0.631±0.063
	PaDiM	0.612±0.080	0.742±0.014	0.787±0.038	0.735±0.031	0.805±0.018	0.792±0.048	0.622±0.013	-	-	-	-
	Patchcore	0.756±0.073	0.923±0.008	0.733±0.002	0.843±0.025	0.885±0.026	0.794±0.040	0.805±0.006	0.833±0.009	0.639±0.010	0.497±0.044	0.739±0.011
	RegAD	0.596±0.040	0.525±0.027	0.596±0.074	0.574±0.042	0.663±0.032	0.571±0.149	0.522±0.050	0.548±0.053	0.534±0.002	0.596±0.075	0.677±0.161
	CoOp	0.781±0.002	0.902±0.006	0.720±0.017	0.818±0.018	0.874±0.017	0.759±0.033	0.860±0.032	0.563±0.004	0.537±0.005	0.618±0.002	0.395±0.008
	WinCLIP	0.754±0.009	0.943±0.004	0.764±0.025	0.858±0.025	0.940±0.021	0.941±0.002	0.912±0.003	0.851±0.010	0.927±0.001	0.632±0.004	0.915±0.003
	Ours (InCTRL)	0.846±0.011	0.975±0.006	0.790±0.018	0.877±0.019	0.945±0.018	0.975±0.016	0.933±0.013	0.902±0.016	0.940±0.010	0.643±0.007	0.928±0.009
	Ours (InCTRL)	0.567±0.034	0.741±0.011	0.708±0.006	0.821±0.042	0.840±0.057	0.794±0.039	0.626±0.022	0.829±0.009	0.849±0.006	0.597±0.028	0.656±0.037
8-shot	SPADE	0.724±0.017	0.769±0.037	0.792±0.025	0.768±0.032	0.820±0.016	0.758±0.025	0.661±0.039	-	-	-	-
	PaDiM	0.837±0.016	0.925±0.003	0.745±0.002	0.860±0.026	0.922±0.019	0.812±0.016	0.817±0.034	0.876±0.004	0.672±0.006	0.526±0.019	0.764±0.004
	Patchcore	0.633±0.027	0.594±0.029	0.603±0.062	0.589±0.040	0.674±0.033	0.632±0.079	0.628±0.026	0.547±0.063	0.555±0.008	0.573±0.076	0.587±0.211
	RegAD	0.817±0.012	0.898±0.005	0.769±0.008	0.822±0.021	0.880±0.014	0.755±0.003	0.914±0.027	0.567±0.007	0.542±0.005	0.619±0.004	0.399±0.006
	CoOp	0.814±0.010	0.941±0.001	0.796±0.015	0.868±0.020	0.947±0.025	0.944±0.001	0.915±0.008	0.867±0.007	0.928±0.001	0.641±0.004	0.916±0.003
	WinCLIP	0.872±0.013	0.978±0.006	0.806±0.036	0.887±0.021	0.953±0.013	0.983±0.012	0.936±0.008	0.920±0.003	0.945±0.002	0.646±0.003	0.934±0.008
	Ours (InCTRL)	0.567±0.034	0.741±0.011	0.708±0.006	0.821±0.042	0.840±0.057	0.794±0.039	0.626±0.022	0.829±0.009	0.849±0.006	0.597±0.028	0.656±0.037
	Ours (InCTRL)	0.724±0.017	0.769±0.037	0.792±0.025	0.768±0.032	0.820±0.016	0.758±0.025	0.661±0.039	-	-	-	-

Table 1. AUROC results(mean±std) on nine real-world AD datasets under various few-shot AD settings. Best results and the second-best results are respectively highlighted in **red** and **blue**. ‘Baseline’ is a WinCLIP-based zero-shot AD model.

$s_p(x)$ is added into Eq. 8 because such patch-level anomaly scores are crucial for detecting local abnormal regions to which the ϕ -based holistic anomaly score can often overlook. α is a hyper-parameter that modulates the contribution of the patch-level residual score. Lastly, we optimize the final anomaly score $s(x)$ using X_{train} :

$$\mathcal{L}_h = \frac{1}{N} \sum_{x \in X_{train}} \mathcal{L}_b(s(x), y_x). \quad (9)$$

Thus, the full InCTRL model is optimized by minimizing the overall loss as follows:

$$\mathcal{L}_{InCTRL} = \mathcal{L}_{IRL} + \mathcal{L}_h. \quad (10)$$

Inference. During inference, for a given test image x_t and the K -shot normal image prompt set \mathcal{P} from the target dataset, they are fed forward through the visual encoder and the adapter layers, obtaining \mathbf{M}_{x_t} and $s_i(x_t)$. The text prompt sets used during training are used to obtain $s_a(x_t)$. Lastly, we obtain the final anomaly score of x_t via Eq. 8.

4. Experiments

4.1. Experimental Setup

Datasets. To verify the efficiency of our method InCTRL, we conduct comprehensive experiments across nine real-world AD datasets, including five industrial defect inspection dataset (MVTec AD [7], VisA [75], AITEX [47], ELPV [18], SDD [48]), two medical image datasets (BrainMRI [43], HeadCT [43]), and two semantic anomaly detection datasets: MNIST [28] and CIFAR-10 [27] under both one-vs-all and multi-class protocols [11, 42]. Under the one-vs-all protocol, one class is used as normal, with

the other classes treated as abnormal; while under the multi-class protocol, images of even-number classes from MNIST and animal-related classes from CIFAR-10 are treated as normal, with the images of the other classes are considered as anomalies (see Appendix A for more details).

To assess the GAD performance, MVTec AD, the combination of its training and test sets, is used as the auxiliary training data, on which GAD models are trained, and they are subsequently evaluated on the test set of the other eight datasets without any further training. We train the model on VisA when evaluating the performance on MVTec AD. The few-shot normal prompts for the target data are randomly sampled from the training set of target datasets and remain the same for all models for fair comparison. We evaluate the performance with the number of few-shot normal prompt set to $K = 2, 4, 8$. The reported results are averaged over three independent runs with different random seeds.

Competing Methods and Evaluation Metrics. Since we aim to achieve a generalist AD model, the comparison is focus on detectors of similar generalist detection capabilities. Following [25], InCTRL is compared with three conventional full-shot AD approaches, including SPADE [16], PaDiM [17], and PatchCore [41], all of which are adapted to the few-shot setting by performing their distance-based anomaly scoring based on the few-shot normal samples. We also compare with state-of-the-art (SotA) conventional few-shot AD method RegAD [23] and the CLIP-driven method WinCLIP [25]. The popular prompt learning method CoOp [72] is used as an additional baseline that is trained on the auxiliary data as InCTRL, after which it uses the few-shot anomaly scoring strategy in WinCLIP to perform anomaly detection.

As for evaluation metrics, following previous works [11, 23, 25, 38, 41], we use two popular metrics AUROC (Area

Under the Receiver Operating Characteristic) and AUPRC (Area Under the Precision-Recall Curve) to evaluate the AD performance. We also evaluate the number of parameters and per-image inference time of CLIP-based methods, which is presented in our Appendix C.1.

Implementation Details. By default, for CLIP-based models, including WinCLIP, CoOp and our InCTRL, we adopt the same CLIP implementation, OpenCLIP [24], and its public pre-trained backbone ViT-B/16+ in our experiments. Adam is used as the optimizer and the initial learning rate is set to 1e-3 by default. The text prompts used in InCTRL are kept exactly the same as WinCLIP. To enable the model to recognize both normal and abnormal objects while preventing overfitting, the training epochs is set to 10 with a batch size of 48 on a single GPU (NVIDIA GeForce RTX 3090). SPADE, PaDiM and WinCLIP¹ use the same image prompts as InCTRL for fair comparison, and the official implementation of PatchCore, RegAD and CoOp is taken. Further details are provided in Appendix B.

4.2. Main Results

Tables 1 and 2 present the comparison results of InCTRL to six SotA competing methods in AUROC and AUPRC, respectively, on nine real-world AD datasets. Note that the results for MVTEC AD, VisA, and the one-vs-all settings of MNIST and CIFAR-10 represent an average result across their respective data subsets (see Appendix C for breakdown results). Below we analyze these results in detail.

Generalization to Industrial Defects. Generally, for the five industrial defect AD datasets, InCTRL significantly outperforms all competing models on almost all cases across the three few-shot settings. With more few-shot image prompts, the performance of all methods generally gets better. Specifically, Patchcore shows better performance than SPADE, PaDiM and RegAD, but all of which generalize badly on these datasets. WinCLIP obtains fairly good generalization and surpasses Patchcore, owing to CLIP’s superior recognition capabilities. Due to the in-context residual information is well transferable across the datasets, InCTRL exhibits superior performance, outperforming the SotA models by a large margin, particularly on challenging datasets like ELPV and SDD. As a result, our InCTRL model respectively gains up to 11.3%, 6.5%, 3.7% AUROC and 6.4%, 5.6%, 6% AUPRC enhancements than the best competing method.

Generalization to Medical Image Anomalies. When applied to medical image AD datasets, InCTRL consistently outperforms SotA models in all few-shot settings. It is evident that all competing methods perform poorly except WinCLIP. Impressively, using only two normal image

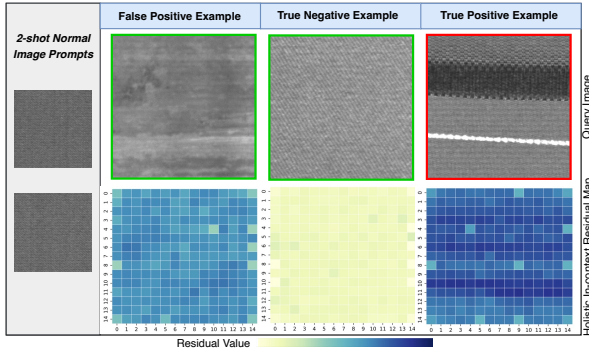


Figure 3. Visualization of query images x_t and their holistic in-context residual maps $M_{x_t}^+$. Green and Red frames indicate normal and abnormal images respectively. Deeper colors in the residual maps represent larger residual values.

prompts, InCTRL can obtain over 97.3% in AUPRC on BrainMRI, despite it does not have any training on medical data. On average, InCTRL surpass the best competing model by 3.9%, 3.4%, 3.9% in AUROC and 0.6%, 1%, 1% in AUPRC for $K = 2, 4, 8$ settings, respectively.

Generalization to Semantic Anomalies under Both One-vs-all and Multi-class Settings.

On detecting semantic anomalies, InCTRL again consistently surpasses all SotA models. Remarkably, InCTRL can obtain 90+% in AUROC when the previous SotA methods can obtain 50%-65% AUROC only, showcasing highly promising GAD performance. Notably, WinCLIP achieves good performance on CIFAR-10. In contrast, CoOp experiences a notable decline, presumably losing crucial semantic knowledge when adapting to the auxiliary data that diverges significantly from the semantic AD task. Overall, our InCTRL achieves the best performance with up to 8.2%, 5.1%, 4.4% AUROC and 2.3%, 1.9%, 2.5% AUPRC improvement compared to the best contender on $K = 2, 4, 8$ settings, respectively.

4.3. Why Does InCTRL Generalize Well?

Ablation Study. We examine the contribution of three key components of InCTRL on the generalization: text prompt-guided features (T), patch-level residuals (P), and image-level residuals (I), as well as their combinations. The results are reported in Table 3. The experiment results indicate that for industrial defect AD datasets, visual residual features play a more significant role compared to text prompt-based features, particularly on datasets like ELPV [18], SDD [48], and AITEX [47]. On the medical image AD datasets, both visual residuals and textual knowledge contribute substantially to performance enhancement, exhibiting a complementary relation. On semantic AD datasets, the results are dominantly influenced by patch-level residuals and/or text prompt-based features. Importantly, our three components are generally mutually complementary, resulting in the superior detection generaliza-

¹No official implementation of WinCLIP is available. Our implementation is available at <https://github.com/mala-lab/WinCLIP>.

Setup	Methods	Industrial Defects					Medical Anomalies		Semantic Anomalies			
		ELPV	SDD	AITEX	VisA	MVTec AD	BrainMRI	HeadCT	One-vs-all		Multi-class	
									MNIST	CIFAR-10	MNIST	CIFAR-10
2-shot	Baseline (0-shot)	0.855±0.000	0.886±0.000	0.552±0.000	0.812±0.000	0.957±0.000	0.988±0.000	0.970±0.000	0.940±0.000	0.990±0.000	0.606±0.000	0.852±0.000
	SPADE	0.618±0.007	0.366±0.105	0.470±0.008	0.818±0.031	0.922±0.023	0.952±0.009	0.851±0.022	0.965±0.004	0.971±0.003	0.615±0.068	0.502±0.035
	PaDiM	0.707±0.058	0.337±0.008	0.529±0.034	0.719±0.027	0.890±0.015	0.902±0.046	0.876±0.017	-	-	-	-
	Patchcore	0.840±0.031	0.676±0.003	0.378±0.008	0.841±0.023	0.939±0.012	0.921±0.017	0.913±0.002	0.956±0.001	0.926±0.002	0.482±0.025	0.574±0.015
	RegAD	0.679±0.005	0.173±0.019	0.275±0.035	0.614±0.037	0.837±0.034	0.872±0.065	0.854±0.009	0.913±0.006	0.909±0.003	0.612±0.013	0.672±0.008
	CoOp	0.841±0.020	0.543±0.004	0.443±0.050	0.835±0.019	0.922±0.007	0.923±0.002	0.937±0.014	0.926±0.003	0.911±0.002	0.607±0.009	0.371±0.013
	WinCLIP	0.849±0.010	0.865±0.004	0.500±0.043	0.859±0.021	0.965±0.007	0.989±0.003	0.975±0.012	0.963±0.001	0.990±0.001	0.614±0.005	0.876±0.016
Ours (InCTRL)	0.913±0.008	0.917±0.009	0.519±0.022	0.877±0.016	0.969±0.004	0.994±0.013	0.981±0.013	0.975±0.004	0.992±0.000	0.618±0.012	0.899±0.010	
4-shot	SPADE	0.627±0.011	0.385±0.018	0.451±0.031	0.826±0.024	0.924±0.015	0.958±0.017	0.854±0.016	0.966±0.008	0.973±0.002	0.611±0.053	0.487±0.047
	PaDiM	0.724±0.067	0.351±0.012	0.540±0.053	0.758±0.018	0.909±0.013	0.956±0.011	0.890±0.011	-	-	-	-
	Patchcore	0.871±0.042	0.703±0.013	0.377±0.001	0.860±0.016	0.950±0.013	0.945±0.017	0.941±0.009	0.972±0.002	0.934±0.003	0.504±0.025	0.606±0.010
	RegAD	0.688±0.018	0.176±0.003	0.294±0.031	0.628±0.034	0.846±0.026	0.900±0.041	0.810±0.028	0.916±0.013	0.908±0.001	0.522±0.085	0.681±0.127
	CoOp	0.867±0.003	0.594±0.014	0.454±0.014	0.842±0.016	0.924±0.008	0.932±0.013	0.957±0.017	0.929±0.002	0.915±0.003	0.611±0.003	0.374±0.012
	WinCLIP	0.864±0.004	0.868±0.003	0.513±0.017	0.875±0.023	0.968±0.008	0.990±0.001	0.974±0.002	0.971±0.002	0.990±0.000	0.611±0.011	0.882±0.009
	Ours (InCTRL)	0.916±0.009	0.924±0.015	0.548±0.016	0.902±0.027	0.972±0.006	0.994±0.013	0.984±0.011	0.980±0.007	0.992±0.004	0.620±0.004	0.901±0.020
8-shot	SPADE	0.641±0.018	0.394±0.024	0.427±0.008	0.844±0.031	0.930±0.016	0.962±0.014	0.860±0.019	0.974±0.002	0.976±0.001	0.613±0.035	0.515±0.024
	PaDiM	0.798±0.014	0.384±0.045	0.555±0.031	0.781±0.024	0.927±0.012	0.946±0.007	0.896±0.009	-	-	-	-
	Patchcore	0.915±0.007	0.708±0.009	0.389±0.003	0.873±0.022	0.962±0.013	0.957±0.007	0.931±0.006	0.979±0.001	0.942±0.002	0.530±0.037	0.635±0.019
	RegAD	0.696±0.015	0.246±0.031	0.314±0.036	0.643±0.032	0.855±0.021	0.908±0.013	0.881±0.014	0.919±0.018	0.911±0.001	0.566±0.048	0.558±0.159
	CoOp	0.905±0.008	0.578±0.001	0.514±0.003	0.848±0.020	0.933±0.007	0.927±0.007	0.965±0.018	0.937±0.004	0.920±0.003	0.610±0.001	0.376±0.003
	WinCLIP	0.897±0.007	0.865±0.001	0.562±0.024	0.880±0.021	0.973±0.009	0.991±0.000	0.975±0.003	0.974±0.001	0.990±0.000	0.616±0.006	0.887±0.006
	Ours (InCTRL)	0.926±0.006	0.925±0.011	0.561±0.034	0.904±0.025	0.977±0.006	0.996±0.003	0.985±0.005	0.989±0.001	0.994±0.001	0.622±0.008	0.912±0.005

Table 2. AUPRC results(mean±std) on nine real-world AD datasets under various few-shot AD settings. Best results and the second-best results are respectively highlighted in **red** and **blue**. ‘Baseline’ is a WinCLIP-based zero-shot AD model.

T	P	I	Industrial Defects				Medical Anomalies		Semantic Anomalies			
			ELPV	SDD	AITEX	VisA	BrainMRI	HeadCT	One-vs-all		Multi-class	
									MNIST	CIFAR-10	MNIST	CIFAR-10
✓	×	×	0.733	0.946	0.733	0.787	0.926	0.900	0.678	0.924	0.620	0.900
×	✓	×	0.794	0.946	0.730	0.843	0.859	0.829	0.890	0.850	0.504	0.894
×	×	✓	0.791	0.931	0.796	0.808	0.898	0.906	0.694	0.704	0.612	0.712
✓	×	×	0.796	0.947	0.711	0.840	0.938	0.919	0.887	0.924	0.622	0.920
✓	×	✓	0.783	0.938	0.756	0.792	0.932	0.924	0.601	0.800	0.627	0.905
×	✓	✓	0.816	0.945	0.785	0.856	0.952	0.918	0.760	0.765	0.618	0.821
✓	✓	✓	0.839	0.972	0.761	0.856	0.973	0.929	0.892	0.935	0.635	0.924
concatenation			0.809	0.951	0.708	0.832	0.955	0.913	0.819	0.863	0.597	0.823
average			0.793	0.940	0.720	0.809	0.941	0.922	0.776	0.881	0.567	0.834

Table 3. AUROC results for ablation study under two-shot setting. Best results and the second-best results are respectively in **red** and **blue**. The results for VisA, and the one-vs-all settings of MNIST and CIFAR-10 represent an average result across their respective data subsets

tion across the datasets.

Significance of In-context Residual Learning. To assess the importance of learning the residuals in InCTRL, we experiment with two alternative operations in both multi-layer patch-level and image-level residual learning: replacing the residual operation with 1) a **concatenation** operation and 2) an **average** operation, with all the other components of InCTRL fixed. As shown in Table 3, the in-context residual learning generalizes much better than the other two alternative ways, significantly enhancing the model’s performance in GAD across three distinct domains.

4.4. Failure Cases

To understand the results of InCTRL better, we provide visualization results illustrating both successful detection and failures by InCTRL on ELPV [18]. As depicted in Figure 3, the incorrectly identified anomalies (False Positive example) shows substantially texture difference compared to the two normal image prompts, leading to similarly large residual values as that for a True Positive example. In contrast, when the query image resembles the normal image prompts well, the residual values are clearly very small, as shown by the True Negative example. These cases may be remedied

when the image prompts include similar normal images as the falsely identified anomaly. This failure case exemplifies the challenge of GAD using only a few image prompts.

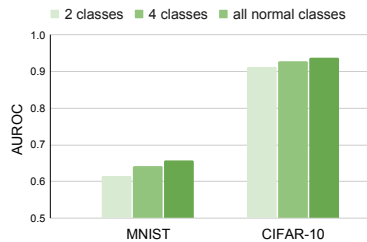


Figure 4. AUROC of InCTRL with sample prompts from varying numbers of normal classes.

This case is related to the diversity of few-shot normal sample prompts, which is particularly important when the normal data is complex, e.g., having multiple different normal patterns in our multi-class protocol.

To rigorously investigate this problem, we evaluate the performance of the two datasets under the multi-class protocol, with varying numbers of normal classes included in the normal image prompts in the eight-shot setting. As shown in Figure 4, the performance of InCTRL continually improves when the prompt sets have samples from more normal classes.

5. Conclusion

In this work we introduce a GAD task to evaluate the generalization capability of AD methods in identifying anomalies across various scenarios without any training on the target datasets. This is the first study dedicated to a generalist approach to anomaly detection, encompassing industrial defects, medical anomalies, and semantic anomalies. Then we propose an approach, called InCTRL, to addressing this problem under a few-shot setting. InCTRL achieves a superior GAD generalization by holistic in-context residual learning. Extensive experiments are performed on nine AD datasets to establish a GAD evaluation benchmark for the aforementioned three popular AD tasks, on which InCTRL significantly and consistently outperforms SotA competing models across multiple few-shot settings.

References

- [1] Abhishek Aich, Kuan-Chuan Peng, and Amit K Roy-Chowdhury. Cross-domain video anomaly detection without target domain adaptation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2579–2591, 2023. 2
- [2] Samet Akcay, Amir Atapour-Abarghouei, and Toby P Breckon. Ganomaly: Semi-supervised anomaly detection via adversarial training. In *Computer Vision—ACCV 2018: 14th Asian Conference on Computer Vision, Perth, Australia, December 2–6, 2018, Revised Selected Papers, Part III 14*, pages 622–637. Springer, 2019. 1, 2
- [3] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022. 3
- [4] Amir Bar, Yossi Gandelsman, Trevor Darrell, Amir Globerson, and Alexei Efros. Visual prompting via image inpainting. *Advances in Neural Information Processing Systems*, 35:25005–25017, 2022. 3
- [5] Niamh Belton, Misgina Tsighe Hagos, Aonghus Lawlor, and Kathleen M Curran. Fewsome: One-class few shot anomaly detection with siamese networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2977–2986, 2023. 2, 3
- [6] Liron Bergman and Yedid Hoshen. Classification-based anomaly detection for general data. *arXiv preprint arXiv:2005.02359*, 2020. 1, 2
- [7] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. Mvtec ad—a comprehensive real-world dataset for unsupervised anomaly detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9592–9600, 2019. 6, 12
- [8] Paul Bergmann, Michael Fauser, David Sattlegger, and Carsten Steger. Uninformed students: Student-teacher anomaly detection with discriminative latent embeddings. In

- Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4183–4192, 2020. 1, 2
- [9] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 3
 - [10] Yu Cai, Hao Chen, Xin Yang, Yu Zhou, and Kwang-Ting Cheng. Dual-distribution discrepancy with self-supervised refinement for anomaly detection in medical images. *Medical Image Analysis*, 86:102794, 2023. 2, 14
 - [11] Tri Cao, Jiawen Zhu, and Guansong Pang. Anomaly detection under distribution shift. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6511–6523, 2023. 1, 2, 6, 12
 - [12] Yunkang Cao, Xiaohao Xu, Jiangning Zhang, Yuqi Cheng, Xiaonan Huang, Guansong Pang, and Weiming Shen. A survey on visual anomaly detection: Challenge, approach, and prospect. *arXiv preprint arXiv:2401.16402*, 2024. 1
 - [13] Ting Chen, Saurabh Saxena, Lala Li, David J Fleet, and Geoffrey Hinton. Pix2seq: A language modeling framework for object detection. *arXiv preprint arXiv:2109.10852*, 2021. 3
 - [14] Ting Chen, Saurabh Saxena, Lala Li, Tsung-Yi Lin, David J Fleet, and Geoffrey E Hinton. A unified sequence interface for vision tasks. *Advances in Neural Information Processing Systems*, 35:31333–31346, 2022. 3
 - [15] Yuanhong Chen, Yu Tian, Guansong Pang, and Gustavo Carneiro. Deep one-class classification via interpolated gaussian descriptor. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 383–392, 2022. 1, 2
 - [16] Niv Cohen and Yedid Hoshen. Sub-image anomaly detection with deep pyramid correspondences. *arXiv preprint arXiv:2005.02357*, 2020. 2, 3, 4, 6
 - [17] Thomas Defard, Aleksandr Setkov, Angélique Loesch, and Romaric Audigier. Padim: a patch distribution modeling framework for anomaly detection and localization. In *International Conference on Pattern Recognition*, pages 475–489. Springer, 2021. 2, 3, 4, 6
 - [18] Sergiu Deitsch, Vincent Christlein, Stephan Berger, Claudia Buerhop-Lutz, Andreas Maier, Florian Gallwitz, and Christian Riess. Automatic classification of defective photovoltaic module cells in electroluminescence images. *Solar Energy*, 185:455–468, 2019. 6, 7, 8, 12
 - [19] Hanqiu Deng and Xingyu Li. Anomaly detection via reverse distillation from one-class embedding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9737–9746, 2022. 1, 2
 - [20] Choubo Ding, Guansong Pang, and Chunhua Shen. Catching both gray and black swans: Open-set supervised anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7388–7398, 2022. 2
 - [21] Yaru Hao, Haoyu Song, Li Dong, Shaohan Huang, Zewen Chi, Wenhui Wang, Shuming Ma, and Furu Wei. Language models are general-purpose interfaces. *arXiv preprint arXiv:2206.06336*, 2022. 3
 - [22] 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8791–8800, 2021. 1, 2
- [23] Chaoqin Huang, Haoyan Guan, Aofan Jiang, Ya Zhang, Michael Spratling, and Yan-Feng Wang. Registration based few-shot anomaly detection. In *European Conference on Computer Vision*, pages 303–319. Springer, 2022. 2, 3, 6
- [24] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, et al. Openclip. *Zenodo*, 4:5, 2021. 7, 13
- [25] Jongheon Jeong, Yang Zou, Taewan Kim, Dongqing Zhang, Avinash Ravichandran, and Onkar Dabeer. Winclip: Zero-/few-shot anomaly classification and segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19606–19616, 2023. 2, 3, 5, 6, 13
- [26] Alexander Kolesnikov, André Susano Pinto, Lucas Beyer, Xiaohua Zhai, Jeremiah Harmsen, and Neil Houlsby. Uvim: A unified modeling approach for vision with learned guiding codes. *Advances in Neural Information Processing Systems*, 35:26295–26308, 2022. 3
- [27] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research). 2009. URL <http://www.cs.toronto.edu/kriz/cifar.html>, 5, 2009. 6, 12, 13
- [28] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 6, 12, 13
- [29] Liu Li, Mai Xu, Xiaofei Wang, Lai Jiang, and Hanruo Liu. Attention based glaucoma detection: A large-scale database and cnn model, 2019. 14
- [30] Jingyi Liao, Xun Xu, Manh Cuong Nguyen, Adam Goodge, and Chuan Sheng Foo. Coft-ad: Contrastive fine-tuning for few-shot anomaly detection. *arXiv preprint arXiv:2402.18998*, 2024. 3
- [31] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 5
- [32] Wenrui Liu, Hong Chang, Bingpeng Ma, Shiguang Shan, and Xilin Chen. Diversity-measurable anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12147–12156, 2023. 1, 2
- [33] Jiasen Lu, Christopher Clark, Rowan Zellers, Roozbeh Motlaghi, and Aniruddha Kembhavi. Unified-io: A unified model for vision, language, and multi-modal tasks. *arXiv preprint arXiv:2206.08916*, 2022. 3
- [34] Yiwei Lu, Frank Yu, Mahesh Kumar Krishna Reddy, and Yang Wang. Few-shot scene-adaptive anomaly detection. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*, pages 125–141. Springer, 2020. 2
- [35] Guansong Pang, Kai Ming Ting, and David Albrecht. Lesinn: Detecting anomalies by identifying least similar nearest neighbours. In *2015 IEEE international conference on data mining workshop (ICDMW)*, pages 623–630. IEEE, 2015. 4
- [36] Guansong Pang, Longbing Cao, Ling Chen, and Huan Liu. Learning representations of ultrahigh-dimensional data for random distance-based outlier detection. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 2041–2050, 2018. 4
- [37] Guansong Pang, Chunhua Shen, Longbing Cao, and Anton Van Den Hengel. Deep learning for anomaly detection: A review. *ACM computing surveys (CSUR)*, 54(2):1–38, 2021. 1
- [38] Guansong Pang, Chunhua Shen, Huidong Jin, and Anton van den Hengel. Deep weakly-supervised anomaly detection. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1795–1807, 2023. 6
- [39] Hyunjong Park, Jongyoun Noh, and Bumsub Ham. Learning memory-guided normality for anomaly detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14372–14381, 2020. 1, 2
- [40] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2, 4
- [41] Karsten Roth, Latha Pemula, Joaquin Zepeda, Bernhard Schölkopf, Thomas Brox, and Peter Gehler. Towards total recall in industrial anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14318–14328, 2022. 2, 3, 4, 6
- [42] Lukas Ruff, Robert A Vandermeulen, Nico Görnitz, Alexander Binder, Emmanuel Müller, Klaus-Robert Müller, and Marius Kloft. Deep semi-supervised anomaly detection. In *ICLR*, 2020. 1, 2, 6
- [43] Mohammadreza Salehi, Niousha Sadjadi, Soroosh Baselizadeh, Mohammad H Rohban, and Hamid R Rabiee. Multiresolution knowledge distillation for anomaly detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14902–14912, 2021. 1, 2, 6, 12, 13
- [44] Thomas Schlegl, Philipp Seeböck, Sebastian M Waldstein, Georg Langs, and Ursula Schmidt-Erfurth. f-anogan: Fast unsupervised anomaly detection with generative adversarial networks. *Medical image analysis*, 54:30–44, 2019. 1, 2
- [45] Eli Schwartz, Assaf Arbelle, Leonid Karlinsky, Sivan Harary, Florian Scheidegger, Sivan Doveh, and Raja Giryes. Maeday: Mae for few and zero shot anomaly-detection. *arXiv preprint arXiv:2211.14307*, 2022. 2, 3
- [46] Shelly Sheynin, Sagie Benaim, and Lior Wolf. A hierarchical transformation-discriminating generative model for few shot anomaly detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8495–8504, 2021. 2, 3
- [47] Javier Silvestre-Blanes, Teresa Albero-Albero, Ignacio Miralles, Rubén Pérez-Llorens, and Jorge Moreno. A public

- fabric database for defect detection methods and results. *Autex Research Journal*, 19(4):363–374, 2019. 6, 7, 12, 13
- [48] 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. 6, 7, 12, 13
- [49] David MJ Tax and Robert PW Duin. Support vector data description. *Machine learning*, 54:45–66, 2004. 1, 2
- [50] Yu Tian, Guansong Pang, Fengbei Liu, Yuanhong Chen, Seon Ho Shin, Johan W Verjans, Rajvinder Singh, and Gustavo Carneiro. Constrained contrastive distribution learning for unsupervised anomaly detection and localisation in medical images. In *Medical Image Computing and Computer Assisted Intervention—MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part V 24*, pages 128–140. Springer, 2021. 2
- [51] Yu Tian, Fengbei Liu, Guansong Pang, Yuanhong Chen, Yuyuan Liu, Johan W Verjans, Rajvinder Singh, and Gustavo Carneiro. Self-supervised pseudo multi-class pre-training for unsupervised anomaly detection and segmentation in medical images. *Medical image analysis*, 90:102930, 2023. 2
- [52] Tran Dinh Tien, Anh Tuan Nguyen, Nguyen Hoang Tran, Ta Duc Huy, Soan Duong, Chanh D Tr Nguyen, and Steven QH Truong. Revisiting reverse distillation for anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24511–24520, 2023. 1, 2
- [53] Guodong Wang, Shumin Han, Errui Ding, and Di Huang. Student-teacher feature pyramid matching for anomaly detection. *arXiv preprint arXiv:2103.04257*, 2021. 1, 2
- [54] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International Conference on Machine Learning*, pages 23318–23340. PMLR, 2022. 3
- [55] Xinlong Wang, Wen Wang, Yue Cao, Chunhua Shen, and Tiejun Huang. Images speak in images: A generalist painter for in-context visual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6830–6839, 2023. 3
- [56] Xinlong Wang, Xiaosong Zhang, Yue Cao, Wen Wang, Chunhua Shen, and Tiejun Huang. Seggpt: Segmenting everything in context. *arXiv preprint arXiv:2304.03284*, 2023. 3
- [57] Ze Wang, Yipin Zhou, Rui Wang, Tsung-Yu Lin, Ashish Shah, and Ser Nam Lim. Few-shot fast-adaptive anomaly detection. *Advances in Neural Information Processing Systems*, 35:4957–4970, 2022. 2, 3
- [58] Jhih-Ciang Wu, Ding-Jie Chen, Chiou-Shann Fuh, and Tyng-Luh Liu. Learning unsupervised metaformer for anomaly detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4369–4378, 2021. 2, 3
- [59] Peng Wu, Xuerong Zhou, Guansong Pang, Yujia Sun, Jing Liu, Peng Wang, and Yanning Zhang. Open-vocabulary video anomaly detection. *arXiv preprint arXiv:2311.07042*, 2023. 2, 3
- [60] Peng Wu, Xuerong Zhou, Guansong Pang, Lingru Zhou, Qingsen Yan, Peng Wang, and Yanning Zhang. Vadclip: Adapting vision-language models for weakly supervised video anomaly detection. *arXiv preprint arXiv:2308.11681*, 2023. 2, 3
- [61] Tiange Xiang, Yixiao Zhang, Yongyi Lu, Alan L Yuille, Chaoyi Zhang, Weidong Cai, and Zongwei Zhou. Squid: Deep feature in-painting for unsupervised anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23890–23901, 2023. 1, 2
- [62] Guoyang Xie, Jingbao Wang, Jiaqi Liu, Feng Zheng, and Yaochu Jin. Pushing the limits of fewshot anomaly detection in industry vision: Graphcore. *arXiv preprint arXiv:2301.12082*, 2023. 2, 3
- [63] Heng Xu, Tianqing Zhu, Lefeng Zhang, Wanlei Zhou, and Philip S Yu. Machine unlearning: A survey. *ACM Computing Surveys*, 56(1):1–36, 2023. 1
- [64] Xudong Yan, Huaidong Zhang, Xuemiao Xu, Xiaowei Hu, and Pheng-Ann Heng. Learning semantic context from normal samples for unsupervised anomaly detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 3110–3118, 2021. 1, 2
- [65] Xincheng Yao, Ruoqi Li, Zefeng Qian, Yan Luo, and Chongyang Zhang. Focus the discrepancy: Intra-and inter-correlation learning for image anomaly detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6803–6813, 2023.
- [66] Xincheng Yao, Chongyang Zhang, Ruoqi Li, Jun Sun, and Zhenyu Liu. One-for-all: Proposal masked cross-class anomaly detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 4792–4800, 2023. 1, 2
- [67] Jihun Yi and Sungroh Yoon. Patch svdd: Patch-level svdd for anomaly detection and segmentation. In *Proceedings of the Asian Conference on Computer Vision*, 2020. 1, 2
- [68] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14183–14193, 2020. 1, 2
- [69] Vitjan Zavrtnik, Matej Kristan, and Danijel Skočaj. Draem—a discriminatively trained reconstruction embedding for surface anomaly detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8330–8339, 2021.
- [70] Vitjan Zavrtnik, Matej Kristan, and Danijel Skočaj. Reconstruction by inpainting for visual anomaly detection. *Pattern Recognition*, 112:107706, 2021. 1, 2
- [71] Xuan Zhang, Shiyu Li, Xi Li, Ping Huang, Jiulong Shan, and Ting Chen. Destseg: Segmentation guided denoising student-teacher for anomaly detection, 2023. 1, 2
- [72] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16816–16825, 2022. 6, 13
- [73] Qihang Zhou, Guansong Pang, Yu Tian, Shibo He, and Jiming Chen. Anomalyclip: Object-agnostic prompt learning for

zero-shot anomaly detection. In *The Twelfth International Conference on Learning Representations*, 2024. 2, 3

- [74] Jiawen Zhu, Choubo Ding, Yu Tian, and Guansong Pang. Anomaly heterogeneity learning for open-set supervised anomaly detection. *arXiv preprint arXiv:2310.12790*, 2023. 2
- [75] Yang Zou, Jongheon Jeong, Latha Pemula, Dongqing Zhang, and Onkar Dabeer. Spot-the-difference self-supervised pre-training for anomaly detection and segmentation. In *European Conference on Computer Vision*, pages 392–408. Springer, 2022. 6, 12, 13

A. Dataset Details

A.1. Data Statistics of Training and Testing

We conduct extensive experiments on nine real-world Anomaly Detection (AD) datasets, including five industrial defect inspection dataset (MVTec AD [7], VisA [75], ELPV [18], SDD [48], AITEX [47]), two medical image datasets (BrainMRI [43], HeadCT [43]), and two semantic anomaly detection datasets: MNIST [28] and CIFAR-10 [27] under both one-vs-all and multi-class protocols [11].

To assess the Genealist Anomaly Detection (GAD) performance, the full dataset of MVTec AD, including both training set and test set, is used as the auxiliary training data, on which AD models are trained, and they are subsequently evaluated on the test set of the other eight datasets without any further training. We train the model on the full dataset of VisA when evaluating the performance on MVTec AD. The few-shot normal prompts for the target data are randomly sampled from the training set of target datasets and remain the same for all models for fair comparison. Table 4 provides the data statistics of MVTec AD and VisA, while Table 5 shows the test set statistics of the rest datasets.

A.2. Industrial Defect Inspection Datasets

MVTec AD [7] is a widely-used dataset that enables researchers to benchmark the performance of anomaly detection methods in the context of industrial inspection applications. The dataset includes over 5,000 images that are divided into 15 object and texture categories. Each category contains a training set of anomaly-free images, as well as a test set that includes images with both defects and defect-free images.

VisA [75] consists of 10,821 high-resolution color images (9,621 normal and 1,200 anomalous samples) covering 12 objects in 3 domains, making it the largest industrial anomaly detection dataset to date. Both image and pixel-level labels are provided. The anomalous images contain various flaws, including surface defects such as scratches, dents, color spots or crack, and structural defects like misplacement or missing parts.

Dataset	Subset	Type	Original Training		Original Test	
			Normal	Anomaly	Normal	Anomaly
MVTec AD	Carpet	Texture	280	28	89	
	Grid	Texture	264	21	57	
	Leather	Texture	245	32	92	
	Tile	Texture	230	33	83	
	Wood	Texture	247	19	60	
	Bottle	Object	209	20	63	
	Capsule	Object	219	23	109	
	Pill	Object	267	26	141	
	Transistor	Object	213	60	40	
	Zipper	Object	240	32	119	
	Cable	Object	224	58	92	
	Hazelnut	Object	391	40	70	
	Metal_nut	Object	220	22	93	
	Screw	Object	320	41	119	
	Toothbrush	Object	60	12	30	
VisA	candle	Object	900	100	100	
	capsules	Object	542	60	100	
	cashew	Object	450	50	100	
	chewinggum	Object	453	50	100	
	fryum	Object	450	50	100	
	macaroni1	Object	900	100	100	
	macaroni2	Object	900	100	100	
	pcb1	Object	904	100	100	
	pcb2	Object	901	100	100	
	pcb3	Object	905	101	100	
	pcb4	Object	904	101	100	
	pipe.fryum	Object	450	50	100	

Table 4. Data statistics of MVTec AD and VisA. When training GAD models with MVTec AD or VisA datasets, we utilize their complete datasets, including both training and test data. In contrast, for testing GAD models, only the test sets of MVTec AD or VisA are employed for inference.

Dataset	Subset	Type	Test set	
			Normal	Anomaly
MNIST	0	Semantical	980	9020
	1	Semantical	1135	8865
	2	Semantical	1,032	8,968
	3	Semantical	1,010	8,990
	4	Semantical	982	9019
	5	Semantical	892	9108
	6	Semantical	958	9042
	7	Semantical	1028	8972
	8	Semantical	974	9026
	9	Semantical	1009	8991
	even_number	Semantical	4926	5074
CIFAR-10	airplane	Semantical	1000	9000
	automobile	Semantical	1000	9000
	bird	Semantical	1000	9000
	cat	Semantical	1000	9000
	deer	Semantical	1000	9000
	dog	Semantical	1000	9000
	frog	Semantical	1000	9000
	horse	Semantical	1000	9000
	ship	Semantical	1000	9000
	truck	Semantical	1000	9000
animal	Semantical	6000	4000	
ELPV	-	Texture	377	715
SDD	-	Texture	286	54
AITEX	-	Texture	564	183
BrainMRI	-	Medical	25	155
HeadCT	-	Medical	25	100

Table 5. Data statistics of seven AD datasets for inference. These datasets are exclusively used for inference purposes, hence only the details of the test sets are provided.

ELPV [18] is a collection of 2,624 high-resolution grayscale images of solar cells extracted from photovoltaic

modules. These images were extracted from 44 different solar modules, and include both intrinsic and extrinsic defects known to reduce the power efficiency of solar modules. In our study, we only use its test set for evaluation.

SDD [48] is a collection of images captured in a controlled industrial environment, using defective production items as the subject. The dataset includes 52 images with visible defects and 347 product images without any defects. In our study, we only use its test set for evaluation.

AITEX [47] is a textile fabric database that comprises 245 images of 7 different fabrics, including 140 defect-free images (20 for each type of fabric) and 105 images with various types of defects. We only use its test set for evaluation.

A.3. Medical Anomaly Detection Datasets

BrainMRI [43] is a dataset for brain tumor detection obtained from magnetic resonance imaging (MRI) of the brain. In our study, we only use its test set for evaluation.

HeadCT [43] is a dataset consisting of 100 normal head CT slices and 100 slices with brain hemorrhage, without distinction between the types of hemorrhage. Each slice is from a different person, providing a diverse set of images for researchers to develop and test algorithms for hemorrhage detection and classification in medical imaging applications. In our study, we only use its test set for evaluation.

A.4. Semantic Anomaly Detection Datasets

MNIST [28] encompasses 70,000 grayscale images of handwritten digits. It serves as a semantic AD dataset in our work, where we utilize its original test set to construct test sets of one-vs-all and multi-class settings. Under the one-vs-all protocol, one of the ten classes is used as normal, with the other classes treated as abnormal; while under the multi-class protocol, images of even-number classes are treated as normal, with the images of the other classes are considered as anomalies. In this case, the category-level label is set to ‘even_number’.

CIFAR-10 [27] consists of 60,000 colour images in 10 classes, with 6,000 images per class. There are 50,000 training images and 10,000 test images. It serves as a semantic AD dataset in our work, where we utilize its original test set to construct test sets of one-vs-all and multi-class settings. Under the one-vs-all protocol, one of the 10 classes is used as normal, with the other classes treated as abnormal; while under the multi-class protocol, images of animal-related classes are treated as normal, with the images of the other classes are considered as anomalies. In this case, the category-level label is set to ‘animal’.

Normal Examples	<i>‘a photo of a flawless [c] for visual inspection.’</i> <i>‘a cropped photo of a perfect [c].’</i> <i>‘a blurry photo of the [c] without defect.’</i> <i>‘a dark photo of the unblemished [c].’</i> <i>‘a jpeg corrupted photo of a [c] without flaw.’</i>
Abnormal Examples	<i>‘a photo of a [c] with flaw for visual inspection.’</i> <i>‘a cropped photo of a [c] with damage.’</i> <i>‘a blurry photo of the [c] with defect.’</i> <i>‘a dark photo of the [c] with flaw.’</i> <i>‘a jpeg corrupted photo of a [c] with defect.’</i>

Table 6. Examples of normal and abnormal text prompts used in I_{nCTRL} and WinCLIP. $[c]$ represents a class label.

B. Implementation Details

B.1. Data Pre-processing

By default, for all CLIP-based models, including WinCLIP [25], CoOp [72], and I_{nCTRL} , we adopt the same CLIP implementation, OpenCLIP [24], and its public pre-trained backbone ViT-B/16+ in our experiments. Our data preprocessing aligns with OpenCLIP across all datasets. Specifically, this involves channel-wise standardization using a predefined mean and standard deviation after scaling RGB images to the range of $[0, 1]$, followed by bicubic resizing based on Pillow library. In addition, we resize the input resolution to 240×240 to match ViT-B/16+. This resizing is also applied to other baseline models for fair comparison, while retaining their original data preprocessing methods (if there are any).

B.2. Network Architectures

In our experiments, the parameters of visual encoder and text encoder of ViT-B/16+ are kept frozen. This model, while being similar in depth to ViT-B/16, increases the dimensions of image embeddings ($768 \rightarrow 896$), text embedding ($512 \rightarrow 640$) and input resolution ($224 \times 224 \rightarrow 240 \times 240$). For the learnable components, to align with ViT-B/16+’s dimensions, the adapter ψ has input and output dimensions set to 896, including a 224-unit hidden layer with ReLU activation. The image-level anomaly classification learner η takes in-context image-level residual features F_x as input and yields a one-dimensional prediction, where η has two hidden layers with 128 and 64 units respectively. The holistic anomaly scoring model ϕ incorporates two hidden layers, projecting a 225-dimensional in-context residual map to generate a final single-dimensional anomaly score.

B.3. Details of Text prompts

The text prompts used in our work are based on the same main text and ensemble strategy to WinCLIP [25] except CIFAR-10 [75]. Table 6 provides several normal and abnormal examples of text prompts used in I_{nCTRL} (see [25] for the full list of the text prompts). The WinCLIP prompts fail to work for natural images like CIFAR-10, so the nor-

Method	Number of Parameters	Inference Time (ms)
CoOp	6,400	197.3±5.6
WinCLIP	0	227.5±0.7
Ours (InCTRL)	334,916	81.7±1.4

Table 7. Number of Parameters and Per-image Inference time.

mal and abnormal text prompts of CIFAR-10 are designed as ‘a photo of [c] for anomaly detection.’ and ‘a photo without [c] for anomaly detection.’, respectively, which are used in both WinCLIP and InCTRL on CIFAR-10. Here, [c] represents a category-level label, e.g., airplane.

B.4. Implementation of Comparison Methods

For the results of competing methods, we re-implemented SPADE, PaDiM, and WinCLIP, while using the official implementations of PatchCore, RegAD, and CoOp. Differing from SPADE’s original $K = 50$ setup, we use $K = 2, 4, 8$ nearest neighbors to match the few-shot setting. For PaDiM, we select the wide_resnet50_2 model, pretrained on ImageNet, as the feature extractor. To ensure fair empirical comparison, we apply the same image prompts as used in InCTRL across all methods. All reported results are the average of three independent runs, each with a different random seed.

C. Detailed Empirical Results

C.1. Complexity of InCTRL vs. Other CLIP-based Methods

The number of parameters and per-image inference time for CLIP-based methods are shown in Table 7.

Our InCTRL and CoOp involve additional training on auxiliary data compared to the training-free method, i.e., WinCLIP. This results in extra trainable parameters during the training phase. However, the extra time consumption leads to significant performance enhancements in InCTRL, and furthermore, the training can be taken offline, so its computation overhead is generally negligible.

Additionally, as Table 7 shows, InCTRL achieves faster inference compared to WinCLIP’s multi-scale few-shot anomaly scoring approach. Although CoOp adopts a similar few-shot anomaly scoring method as WinCLIP, it gains better efficiency by avoiding the ensemble text prompt strategy in WinCLIP. Result indicates the effectiveness of the InCTRL framework in enhancing the base model’s generalization ability.

C.2. Comparison with Traditional Medical Anomaly Detection Methods

Even though our comparison is focused on detectors of similar generalist detection capabilities, we compare five recent AD methods specifically designed for medical im-

	Dataset	Brain Tumor MRI		LAG	
		AUR	PR	AUR	PR
Traditional Medical AD	FPI	0.831	0.789	0.543	0.556
	PII	0.843	0.805	0.610	0.607
	F-AnoGAN	0.825	0.743	0.842	0.775
	AEU	0.940	0.890	0.813	0.789
	AEU+DDAD	0.942	0.919	0.860	0.840
Generalist AD	WinCLIP	0.779	0.878	0.571	0.731
	Ours (InCTRL)	0.951	0.968	0.832	0.880

Table 8. Comparison with Traditional Medical AD Methods.

ages on two other medical datasets (Brain Tumor MRI² and LAG [29]) in Table 8. The results of traditional medical anomaly detection methods are from Cai *et al.* [10] based on **full-shot** one-class classification setting. It should be noted that better performance is gained if **extra anomaly image data** is used, but it would be very unfair comparison to our method that uses only **8-shot** normal images.

As shown in Table 8, our method outperforms all competing models on almost all cases, despite the fact that our method uses only 8-shot normal images for prompting and does not require any training on the medical data whereas the medical image AD methods require extensive training on a large set of normal medical images, indicating the superior generalized AD capability of our model.

C.3. Full Results on VisA and MVTEC AD

Table 9 presents detailed comparison results of InCTRL against six SotA methods across each category of the VisA dataset. Overall, InCTRL markedly surpasses all competitors in every case within the three few-shot settings. We observe a general improvement in performance across all methods with an increase in the number of few-shot image prompts.

Similarly, Table 10 details the results of InCTRL and six SotA methods across each category of the MVTEC AD dataset. InCTRL again consistently outperforms all baseline models in all few-shot settings.

C.4. Full Results on One-vs-all Setting of Sematic Anomaly Detection Datasets

Table 11 provides the detailed empirical comparison results of InCTRL against five SotA methods using the one-vs-all protocol on the MNIST dataset. Similarly, Table 12 details the performance of InCTRL relative to five SotA methods under the one-vs-all protocol on the CIFAR-10 dataset. InCTRL consistently surpasses all baseline methods in all few-shot settings on both datasets.

²The dataset is available at <https://www.kaggle.com/datasets/masoudnickparvar/brain-tumor-mri-dataset>.

Data subset		AUROC							AUPRC						
		SPADE	PaDiM	Patchcore	RegAD	CoOp	WinCLIP	Ours (InCTRL)	SPADE	PaDiM	Patchcore	RegAD	CoOp	WinCLIP	Ours (InCTRL)
2-shot	candle	0.840±0.057	0.876±0.017	0.935±0.013	0.455±0.042	0.943±0.019	0.974±0.007	0.916±0.006	0.889±0.033	0.839±0.015	0.940±0.021	0.461±0.038	0.950±0.015	0.974±0.006	0.920±0.008
	capsules	0.702±0.106	0.540±0.036	0.573±0.029	0.420±0.030	0.697±0.021	0.772±0.015	0.820±0.029	0.771±0.067	0.664±0.022	0.698±0.016	0.583±0.022	0.822±0.017	0.789±0.011	0.846±0.017
	cashew	0.952±0.033	0.674±0.024	0.915±0.010	0.811±0.021	0.933±0.008	0.943±0.008	0.978±0.007	0.954±0.021	0.814±0.018	0.959±0.005	0.879±0.019	0.971±0.004	0.974±0.006	0.990±0.007
	chewinggum	0.886±0.021	0.734±0.020	0.957±0.008	0.511±0.013	0.973±0.011	0.988±0.009	0.997±0.010	0.909±0.022	0.839±0.015	0.983±0.008	0.703±0.011	0.987±0.008	0.994±0.004	0.998±0.012
	fryum	0.846±0.025	0.678±0.039	0.811±0.024	0.592±0.035	0.899±0.008	0.791±0.037	0.953±0.011	0.912±0.019	0.801±0.021	0.853±0.017	0.743±0.026	0.957±0.005	0.894±0.019	0.976±0.009
	macaroni1	0.690±0.043	0.538±0.052	0.763±0.055	0.425±0.119	0.816±0.025	0.885±0.022	0.770±0.038	0.651±0.038	0.540±0.038	0.795±0.039	0.433±0.082	0.812±0.015	0.893±0.018	0.786±0.013
	macaroni2	0.586±0.024	0.614±0.022	0.560±0.054	0.476±0.023	0.774±0.047	0.663±0.022	0.726±0.023	0.548±0.020	0.606±0.039	0.597±0.031	0.502±0.021	0.774±0.034	0.678±0.020	0.734±0.015
	pcb1	0.883±0.049	0.667±0.028	0.769±0.074	0.653±0.025	0.521±0.032	0.831±0.046	0.937±0.039	0.863±0.043	0.647±0.018	0.684±0.059	0.656±0.016	0.582±0.024	0.831±0.032	0.945±0.024
	pcb2	0.783±0.065	0.600±0.031	0.822±0.012	0.453±0.082	0.673±0.067	0.655±0.025	0.660±0.020	0.789±0.031	0.606±0.021	0.848±0.009	0.476±0.065	0.683±0.033	0.664±0.037	0.669±0.018
	pcb3	0.761±0.034	0.709±0.029	0.803±0.016	0.463±0.027	0.709±0.022	0.759±0.016	0.734±0.018	0.776±0.023	0.656±0.018	0.821±0.012	0.523±0.023	0.728±0.020	0.771±0.021	0.716±0.030
	pcb4	0.885±0.079	0.671±0.087	0.952±0.021	0.763±0.038	0.778±0.043	0.895±0.063	0.966±0.027	0.893±0.057	0.672±0.031	0.955±0.016	0.756±0.025	0.773±0.034	0.891±0.057	0.965±0.021
pipe_fryum	0.727±0.028	0.864±0.024	0.945±0.011	0.667±0.020	0.958±0.008	0.942±0.014	0.835±0.012	0.855±0.017	0.941±0.015	0.964±0.006	0.657±0.019	0.979±0.009	0.960±0.005	0.975±0.006	
MEAN	0.795±0.045	0.680±0.042	0.817±0.028	0.557±0.063	0.806±0.023	0.842±0.024	0.858±0.022	0.818±0.031	0.719±0.027	0.841±0.023	0.614±0.037	0.835±0.019	0.859±0.021	0.877±0.016	
4-shot	candle	0.885±0.043	0.885±0.028	0.947±0.011	0.468±0.037	0.950±0.008	0.975±0.005	0.937±0.005	0.894±0.034	0.841±0.019	0.947±0.007	0.478±0.024	0.955±0.012	0.975±0.007	0.945±0.004
	capsules	0.711±0.091	0.549±0.044	0.716±0.023	0.446±0.022	0.766±0.018	0.813±0.022	0.859±0.005	0.772±0.077	0.666±0.030	0.805±0.014	0.582±0.027	0.863±0.010	0.828±0.013	0.864±0.010
	cashew	0.970±0.018	0.691±0.016	0.952±0.013	0.817±0.015	0.945±0.013	0.955±0.008	0.978±0.013	0.975±0.014	0.832±0.019	0.976±0.006	0.884±0.016	0.978±0.011	0.984±0.005	0.990±0.005
	chewinggum	0.909±0.024	0.838±0.027	0.973±0.006	0.535±0.011	0.975±0.005	0.991±0.007	0.998±0.006	0.911±0.018	0.898±0.029	0.989±0.003	0.716±0.013	0.988±0.006	0.996±0.003	0.998±0.007
	fryum	0.853±0.022	0.725±0.018	0.830±0.019	0.593±0.028	0.898±0.006	0.818±0.027	0.967±0.010	0.914±0.013	0.838±0.013	0.927±0.005	0.744±0.023	0.957±0.006	0.904±0.008	0.985±0.007
	macaroni1	0.714±0.037	0.631±0.043	0.767±0.037	0.456±0.096	0.811±0.019	0.897±0.022	0.776±0.031	0.672±0.025	0.602±0.024	0.789±0.023	0.459±0.048	0.809±0.012	0.901±0.007	0.812±0.012
	macaroni2	0.592±0.012	0.668±0.015	0.548±0.020	0.521±0.034	0.775±0.038	0.672±0.017	0.746±0.018	0.558±0.016	0.646±0.017	0.528±0.016	0.542±0.022	0.774±0.024	0.685±0.026	0.772±0.018
	pcb1	0.892±0.037	0.779±0.021	0.759±0.048	0.693±0.021	0.527±0.024	0.858±0.024	0.959±0.024	0.873±0.032	0.743±0.017	0.672±0.021	0.694±0.015	0.585±0.015	0.875±0.027	0.962±0.021
	pcb2	0.798±0.055	0.608±0.024	0.872±0.011	0.455±0.068	0.689±0.045	0.678±0.016	0.699±0.014	0.790±0.046	0.605±0.025	0.886±0.013	0.482±0.060	0.691±0.028	0.690±0.035	0.724±0.013
	pcb3	0.770±0.031	0.709±0.012	0.825±0.016	0.473±0.013	0.715±0.016	0.779±0.009	0.761±0.016	0.780±0.024	0.657±0.011	0.847±0.017	0.540±0.021	0.735±0.014	0.782±0.019	0.795±0.023
	pcb4	0.889±0.048	0.857±0.043	0.989±0.013	0.765±0.022	0.801±0.033	0.905±0.071	0.975±0.022	0.902±0.013	0.815±0.032	0.988±0.005	0.765±0.023	0.783±0.017	0.907±0.041	0.982±0.028
pipe_fryum	0.745±0.025	0.885±0.019	0.937±0.006	0.668±0.012	0.962±0.010	0.958±0.011	0.869±0.011	0.866±0.017	0.947±0.015	0.971±0.007	0.653±0.018	0.982±0.032	0.968±0.006	0.989±0.005	
MEAN	0.811±0.040	0.735±0.031	0.843±0.025	0.574±0.042	0.818±0.018	0.858±0.025	0.877±0.019	0.826±0.024	0.758±0.018	0.860±0.016	0.628±0.034	0.842±0.016	0.875±0.023	0.902±0.027	
8-shot	candle	0.899±0.036	0.903±0.022	0.966±0.009	0.781±0.031	0.960±0.003	0.980±0.007	0.955±0.006	0.925±0.012	0.858±0.016	0.968±0.004	0.744±0.026	0.964±0.013	0.980±0.008	0.956±0.005
	capsules	0.712±0.104	0.560±0.046	0.756±0.021	0.469±0.019	0.781±0.018	0.815±0.020	0.863±0.025	0.788±0.067	0.691±0.021	0.845±0.012	0.624±0.018	0.877±0.010	0.832±0.011	0.869±0.013
	cashew	0.977±0.026	0.787±0.023	0.958±0.017	0.519±0.013	0.934±0.015	0.961±0.004	0.979±0.015	0.983±0.010	0.882±0.017	0.978±0.008	0.681±0.015	0.972±0.007	0.987±0.006	0.990±0.005
	chewinggum	0.912±0.017	0.890±0.021	0.980±0.009	0.556±0.009	0.973±0.007	0.994±0.008	0.998±0.006	0.926±0.008	0.937±0.016	0.992±0.003	0.720±0.010	0.987±0.005	0.997±0.004	0.999±0.003
	fryum	0.857±0.023	0.767±0.015	0.858±0.017	0.743±0.022	0.901±0.006	0.826±0.023	0.970±0.010	0.919±0.016	0.848±0.015	0.938±0.005	0.846±0.020	0.958±0.007	0.916±0.018	0.990±0.006
	macaroni1	0.728±0.036	0.645±0.038	0.798±0.036	0.463±0.085	0.811±0.012	0.905±0.032	0.751±0.017	0.705±0.026	0.621±0.023	0.799±0.014	0.466±0.056	0.808±0.023	0.913±0.016	0.818±0.018
	macaroni2	0.598±0.018	0.683±0.013	0.547±0.024	0.543±0.039	0.775±0.035	0.677±0.015	0.715±0.017	0.578±0.023	0.651±0.013	0.549±0.012	0.527±0.018	0.775±0.029	0.689±0.026	0.773±0.008
	pcb1	0.897±0.032	0.812±0.016	0.753±0.045	0.531±0.023	0.540±0.025	0.904±0.024	0.964±0.022	0.886±0.026	0.767±0.019	0.659±0.018	0.548±0.022	0.610±0.028	0.877±0.013	0.967±0.004
	pcb2	0.804±0.047	0.638±0.019	0.907±0.011	0.509±0.079	0.695±0.042	0.690±0.013	0.696±0.019	0.815±0.032	0.613±0.026	0.912±0.013	0.563±0.065	0.696±0.034	0.696±0.026	0.727±0.017
	pcb3	0.777±0.028	0.749±0.021	0.839±0.019	0.466±0.019	0.716±0.017	0.789±0.006	0.798±0.029	0.795±0.021	0.718±0.012	0.855±0.009	0.478±0.015	0.733±0.012	0.791±0.013	0.796±0.028
	pcb4	0.894±0.045	0.878±0.039	0.993±0.011	0.703±0.023	0.813±0.029	0.918±0.160	0.977±0.028	0.918±0.012	0.838±0.033	0.993±0.004	0.634±0.021	0.808±0.033	0.918±0.009	0.975±0.014
pipe_fryum	0.801±0.026	0.899±0.014	0.961±0.009	0.786±0.014	0.967±0.015	0.959±0.009	0.909±0.018	0.887±0.018	0.953±0.010	0.982±0.008	0.883±0.013	0.987±0.016	0.968±0.004	0.990±0.009	
MEAN	0.821±0.042	0.768±0.032	0.860±0.026	0.589±0.040	0.822±0.021	0.868±0.020	0.887±0.021	0.844±0.031	0.781±0.024	0.873±0.022	0.643±0.032	0.848±0.020	0.880±0.021	0.904±0.025	

Table 9. Fine-grained AUROC and AUPRC results(mean±std) on VisA datasets under various few-shot AD settings. Best results and the second-best results are respectively highlighted in red and blue.

Data subset		AUROC							AUPRC						
		SPADE	PaDiM	Patchcore	RegAD	CoOp	WinCLIP	Ours (InCTRL)	SPADE	PaDiM	Patchcore	RegAD	CoOp	WinCLIP	Ours (InCTRL)
2-shot	carpet	0.976±0.008	0.992±0.005	0.996±0.003	0.601±0.051	1.000±0.000	0.998±0.002	0.997±0.003	0.962±0.005	0.998±0.003	0.999±0.001	0.853±0.023	1.000±0.000	1.000±0.000	0.999±0.001
	grid	0.404±0.054	0.560±0.041	0.574±0.102	0.627±0.069	0.675±0.052	0.965±0.007	0.974±0.005	0.727±0.036	0.751±0.038	0.785±0.037	0.830±0.029	0.875±0.038	0.984±0.004	0.988±0.003
	leather	1.000±0.000	1.000±0.000	1.000±0.000	0.556±0.047	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000	0.835±0.016	1.000±0.000	1.000±0.000	1.000±0.000
	tile	0.940±0.011	0.947±0.013	0.987±0.012	0.807±0.027	0.999±0.002	0.997±0.003	1.000±0.000	0.978±0.007	0.980±0.005	0.996±0.004	0.925±0.013	0.998±0.002	0.999±0.002	1.000±0.000
	wood	0.932±0.008	0.982±0.005	0.985±0.009	0.921±0.011	0.966±0.008	0.993±0.005	0.996±0.003	0.979±0.004	0.995±0.002	0.995±0.005	0.976±0.014	0.989±0.006	0.998±0.003	0.998±0.002
	bottle	1.000±0.000	0.971±0.012	1.000±0.000	0.529±0.045	0.973±0.013	0.993±0.007	0.998±0.002	1.000±0.000	0.991±0.005	1.000±0.000	0.830±0.026	0.992±0.004	0.998±0.003	1.000±0.000
	capsule	0.736±0.027	0.741±0.026	0.666±0.068	0.593±0.073	0.823±0.024	0.806±0.068	0.860±0.081	0.935±0.015	0.945±0.010	0.911±0.021	0.861±0.017	0.959±0.016	0.945±0.042	0.960±0.063
	pill	0.759±0.023	0.601±0.059	0.829±0.022	0.693±0.036	0.895±0.022	0.888±0.014	0.893±0.024	0.944±0.018	0.889±0.031	0.965±0.009	0.923±0.015	0.980±0.009	0.978±0.006	0.973±0.011
	transistor	0.853±0.019	0.791±0.043	0.992±0.014	0.687±0.033	0.702±0.013	0.898±0.021	0.875±0.027	0.845±0.021	0.700±0.055	0.989±0.007	0.694±0.029	0.664±0.025	0.858±0.017	0.845±0.023
	zipper	0.905±0.033	0.767±0.024	0.978±0.016	0.602±0.049	0.877±0.018	0.955±0.013	0.954±0.010	0.974±0.010	0.875±0.021	0.994±0.005	0.869±0.018	0.968±0.011	0.979±0.009	0.988±0.009
	cable	0.767±0.046	0.525±0.057	0.805±0.025	0.513±0.057	0.724±0.067	0.862±0.011	0.884±0.013	0.968±0.015	0.693±0.063	0.893±0.018	0.674±0.023	0.777±0.052	0.904±0.009	0.920±0.009
	hazelnut	0.900±0.035	0.947±0.012	0.864±0.027	0.697±0.035	0.938±0.014	0.950±0.008	0.968±0.006	0.946±0.013	0.964±0.009	0.919±0.015	0.844±0.026	0.959±0.010	0.975±0.005	0.987±0.006
	metal_nut	0.760±0.062	0.613±0.035	0.851±0.053	0.632±0.038	0.942±0.009	0.930±0.011	0.942±0.028	0.876±0.042	0.876±0.016	0.965±0.008	0.878±0.017	0.987±0.006	0.984±0.006	0.985±0.010
	screw	0.531±0.024	0.501±0.046	0.531±0.027	0.654±0.041	0.659±0.027	0.780±0.019	0.803±0.043	0.766±0.011	0.751±0.017	0.747±0.033	0.843±0.022	0.867±0.018	0.903±0.013	0.917±0.022
toothbrush	0.786±0.041	0.844±0.043	0.811±0.029	0.486±0.053	0.703±0.022	0.952±0.021	0.960±0.028	0.937±0.020	0.938±0.011	0.922±0.016	0.717±0.024	0.811±0.019	0.972±0.007	0.975±0.019	
MEAN	0.817±0.054	0.785±0.025	0.858±0.034	0.640±0.047	0.888±0.016	0.931±0.019	0.940±0.015	0.922±0.023	0.890±0.015	0.939±0.012	0.837±0.034	0.922±0.007	0.965±0.007	0.969±0.004	
4-shot	carpet	0.981±0.006	0.994±0.004	0.994±0.005	0.604±0.046	0.998±0.002	0.997±0.003	1.000±0.000	0.996±0.005	0.998±0.002	0.998±0.002	0.855±0.025	0.999±0.001	0.999±0.001	1.000±0.000
	grid	0.476±0.046	0.556±0.049	0.565±0.077	0.660±0.043	0.769±0.035	0.978±0.004	0.998±0.003	0.714±0.038	0.752±0.042	0.781±0.032	0.859±0.025	0.894±0.021	0.985±0.004	0.999±0.001
	leather	1.000±0.000	1.000±0.000	1.000±0.000	0.593±0.037	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000	0.879±0.021	1.000±0.000	1.000±0.000	1.000±0.000
	tile	1.000±0.000	0.954±0.008	0.987±0.010	0.708±0.033	0.999±0.001	0.999±0.001	0.999±0.002	1.000±0.000	0.984±0.004	0.996±0.005	0.883±0.024	1.000±0.000	0.999±0.001	0.999±0.001
	wood	0.989±0.011	0.978±0.006	0.987±0.008	0.887±0.026	0.970±0.006	0.997±0.002	0.995±0.005	0.995±0.006	0.994±0.004	0.996±0.006	0.965±0.018	0.991±0.007	0.998±0.001	0.994±0.004
	bottle	1.000±0.000	0.976±0.008	1.000±0.000	0.656±0.035	0.976±0.011	0.992±0.006	1.000±0.000	1.000±0.000	0.993±0.005	1.000±0.000	0.865±0.031	0.993±0.005	0.997±0.004	1.000±0.000
	capsule	0.729±0.018	0.752±0.019	0.777±0.055	0.610±0.063	0.828±0.019	0.840±0.052	0.868±0.043	0.936±0.013	0.958±0.013	0.943±0.023	0.865±0.014	0.959±0.014	0.948±0.037	0.965±0.027
	pill	0.872±0.021	0.650±0.047	0.875±0.019	0.691±0.028	0.906±0.016	0.890±0.017	0.874±0.019	0.978±0.015	0.926±0.025	0.976±0.010	0.926±0.013	0.982±0.011	0.978±0.005	0.967±0.016
	transistor	0.855±0.016	0.868±0.033	0.994±0.008	0.684±0.025	0.719±0.008	0.898±0.014	0.895±0.021	0.844±0.014	0.812±0.036	0.992±0.006	0.754±0.021	0.671±0.018	0.858±0.015	0.875±0.016
	zipper	0.913±0.027	0.759±0.017	0.979±0.013	0.551±0.027	0.881±0.021	0.965±0.010	0.964±0.008	0.976±0.008	0.872±0.015	0.994±0.005	0.854±0.016	0.961±0.016	0.980±0.007	0.990±0.006
	cable	0.818±0.043	0.550±0.042	0.822±0.032	0.510±0.044	0.751±0.046	0.865±0.008	0.888±0.021	0.904±0.026	0.732±0.048	0.908±0.022	0.663±0.018	0.781±0.034	0.912±0.006	0.920±0.006
	hazelnut	0.771±0.037	0.951±0.030	0.999±0.002	0.685±0.029	0.971±0.010	0.954±0.008	0.972±0.005	0.858±0.034	0.978±0.006	0.999±0.001	0.876±0.020	0.975±0.012	0.987±0.005	0.989±0.006
	metal_nut	0.734±0.048	0.741±0.021	0.944±0.036	0.701±0.032	0.943±0.006	0.948±0.009	0.946±0.024	0.971±0.033	0.923±0.013	0.987±0.004	0.880±0.015	0.987±0.005	0.986±0.005	0.989±0.008
	screw	0.531±0.019	0.504±0.038	0.537±0.026	0.599±0.034	0.720±0.018	0.800±0.013	0.815±0.031	0.760±0.008	0.760±0.014	0.751±0.028	0.806±0.026	0.864±0.014	0.915±0.009	0.920±0.016
toothbrush	0.750±0.042	0.847±0.035	0.818±0.023	0.806±0.041	0.678±0.016	0.975±0.016	0.963±0.024	0.933±0.021	0.949±0.009	0.925±0.013	0.757±0.022	0.804±0.015	0.985±0.005	0.978±0.015	
MEAN	0.828±0.044	0.805±0.018	0.885±0.026	0.663±0.032	0.874±0.017	0.940±0.021	0.945±0.018	0.924±0.015	0.909±0.013	0.950±0.013	0.846±0.026	0.924±0.008	0.968±0.008	0.972±0.006	
8-shot	carpet	0.987±0.004	0.994±0.004	0.993±0.006	0.601±0.045	1.000±0.000	0.998±0.002	1.000±0.000	0.991±0.006	0.998±0.002	0.998±0.002	0.858±0.032	1.000±0.000	0.999±0.001	1.000±0.000
	grid	0.477±0.044	0.612±0.042	0.652±0.053	0.679±0.036	0.759±0.033	0.986±0.006	0.999±0.001	0.718±0.035	0.808±0.033	0.828±0.027	0.862±0.024	0.891±0.017	0.987±0.005	0.999±0.001
	leather	1.000±0.000	1.000±0.000	1.000±0.000	0.646±0.028	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000	1.000±0.000	0.821±0.014	1.000±0.000	1.000±0.000	1.000±0.000
	tile	1.000±0.000	0.954±0.006	0.989±0.006	0.782±0.025	0.999±0.001	0.998±0.001	0.998±0.001	1.000±0.000	0.985±0.003	0.997±0.004	0.952±0.020	1.000±0.000	0.999±0.001	0.999±0.001
	wood	0.991±0.006	0.989±0.004	0.991±0.005	0.909±0.019	0.968±0.007	0.996±0.003	0.996±0.004	0.995±0.005	0.997±0.003	0.997±0.004	0.972±0.013	0.990±0.006	0.999±0.001	0.995±0.003
	bottle	1.000±0.000	0.987±0.007	1.000±0.000	0.659±0.027	0.976±0.010	0.994±0.003	1.000±0.000	1.000±0.000	0.996±0.004	1.000±0.000	0.868±0.025	0.993±0.006	0.998±0.002	1.000±0.000
	capsule	0.729±0.013	0.780±0.015	0.929±0.038	0.689±0.049	0.836±0.015	0.884±0.041	0.890±0.037	0.936±0.012	0.979±0.011	0.983±0.016	0.883±0.011	0.961±0.008	0.968±0.035	0.966±0.022
	pill	0.883±0.029	0.657±0.039	0.891±0.015	0.675±0.024	0.900±0.011	0.895±0.016	0.888±0.016	0.983±0.013	0.939±0.027	0.978±0.011	0.918±0.014	0.981±0.009	0.978±0.007	0.969±0.015
	transistor	0.889±0.026	0.887±0.041	0.997±0.005	0.689±0.018	0.742±0.019	0.898±0.013	0.903±0.019	0.877±0.019	0.864±0.047	0.995±0.005	0.784±0.016	0.712±0.021	0.861±0.013	0.884±0.013
	zipper	0.924±0.022	0.757±0.015	0.984±0.008	0.581±0.021	0.889±0.016	0.965±0.008	0.970±0.012	0.984±0.006	0.880±0.013	0.996±0.004	0.864±0.013	0.965±0.011	0.981±0.005	0.993±0.008
	cable	0.872±0.067	0.602±0.033	0.944±0.027	0.521±0.035	0.819±0.035	0.886±0.007	0.932±0.015	0.931±0.035	0.798±0.037	0.969±0.015	0.675±0.016	0.863±0.027	0.928±0.005	0.960±0.007
	hazelnut	0.784±0.033	0.953±0.026	1.000±0.000	0.744±0.024	0.968±0.012	0.958±0.007	0.975±0.007	0.869±0.030	0.985±0.007	1.000±0.000	0.889±0.021	0.971±0.010	0.980±0.009	0.987±0.008
	metal_nut	0.764±0.036	0.783±0.025	0.969±0.023	0.735±0.035	0.946±0.005	0.948±0.012	0.947±0.026	0.973±0.030	0.955±0.018	0.993±0.004	0.898±0.018	0.988±0.008	0.988±0.006	0.990±0.005
	screw	0.544±0.021	0.501±0.034	0.645±0.021	0.600±0.033	0.718±0.014	0.815±0.014	0.825±0.029	0.764±0.013	0.765±0.016	0.762±0.024	0.809±0.024	0.860±0.012	0.934±0.011	0.935±0.013
toothbrush	0.756±0.036	0.849±0.027	0.850±0.016	0.603±0.075	0.681±0.018	0.978±0.012	0.975±0.022	0.934±0.019	0.956±0.008	0.941±0.011	0.766±0.041	0.818±0.016	0.990±0.003	0.982±0.013	
MEAN	0.840±0.057	0.820±0.016	0.922±0.019	0.674±0.033	0.880±0.014	0.947±0.025									

Data subset		AUROC						AUPRC					
		SPADE	Patchcore	RegAD	CoOp	WinCLIP	Ours (InCTRL)	SPADE	Patchcore	RegAD	CoOp	WinCLIP	Ours (InCTRL)
2-shot	0	0.772±0.174	0.803±0.028	0.615±0.037	0.508±0.031	0.801±0.026	0.970±0.004	0.959±0.036	0.950±0.003	0.923±0.013	0.927±0.004	0.954±0.004	0.995±0.000
	1	0.961±0.021	0.988±0.006	0.734±0.024	0.439±0.018	0.928±0.011	0.868±0.047	0.994±0.004	0.998±0.001	0.956±0.007	0.901±0.002	0.986±0.003	0.935±0.019
	2	0.670±0.055	0.643±0.035	0.658±0.062	0.570±0.053	0.723±0.016	0.860±0.035	0.946±0.010	0.933±0.008	0.935±0.016	0.923±0.005	0.943±0.004	0.969±0.005
	3	0.808±0.065	0.812±0.015	0.716±0.077	0.767±0.042	0.786±0.091	0.954±0.025	0.973±0.010	0.972±0.001	0.952±0.012	0.965±0.014	0.958±0.019	0.992±0.004
	4	0.747±0.036	0.774±0.024	0.319±0.023	0.599±0.026	0.756±0.086	0.965±0.026	0.959±0.007	0.962±0.004	0.870±0.020	0.939±0.016	0.952±0.017	0.995±0.003
	5	0.759±0.036	0.747±0.024	0.427±0.059	0.745±0.027	0.824±0.023	0.914±0.048	0.969±0.005	0.959±0.003	0.894±0.008	0.962±0.007	0.970±0.004	0.988±0.006
	6	0.725±0.045	0.611±0.047	0.341±0.036	0.428±0.042	0.878±0.018	0.779±0.027	0.959±0.005	0.934±0.009	0.873±0.008	0.902±0.005	0.982±0.002	0.952±0.004
	7	0.708±0.123	0.795±0.039	0.428±0.142	0.651±0.023	0.827±0.019	0.962±0.011	0.947±0.027	0.965±0.007	0.887±0.022	0.945±0.003	0.963±0.004	0.995±0.001
	8	0.800±0.102	0.629±0.045	0.588±0.078	0.344±0.025	0.712±0.031	0.702±0.136	0.970±0.016	0.928±0.011	0.936±0.015	0.873±0.005	0.948±0.005	0.943±0.043
	9	0.841±0.020	0.758±0.041	0.427±0.048	0.515±0.036	0.860±0.058	0.950±0.019	0.976±0.003	0.962±0.008	0.900±0.006	0.919±0.013	0.975±0.011	0.989±0.007
	MEAN	0.779±0.024	0.756±0.004	0.525±0.030	0.557±0.006	0.810±0.008	0.892±0.009	0.965±0.004	0.956±0.001	0.913±0.006	0.926±0.003	0.963±0.001	0.975±0.004
4-shot	0	0.854±0.049	0.874±0.003	0.624±0.045	0.509±0.018	0.813±0.014	0.974±0.003	0.978±0.007	0.982±0.001	0.926±0.009	0.928±0.003	0.956±0.003	0.995±0.001
	1	0.955±0.032	0.995±0.002	0.740±0.023	0.446±0.015	0.942±0.016	0.902±0.017	0.993±0.005	0.999±0.000	0.958±0.011	0.913±0.003	0.989±0.003	0.953±0.009
	2	0.658±0.044	0.712±0.044	0.690±0.058	0.581±0.035	0.752±0.029	0.869±0.071	0.944±0.008	0.946±0.011	0.941±0.007	0.927±0.005	0.949±0.005	0.962±0.015
	3	0.837±0.014	0.869±0.026	0.730±0.106	0.771±0.037	0.929±0.047	0.959±0.017	0.978±0.002	0.980±0.005	0.922±0.014	0.967±0.011	0.987±0.009	0.994±0.005
	4	0.790±0.033	0.856±0.051	0.335±0.034	0.608±0.022	0.841±0.017	0.972±0.022	0.919±0.081	0.976±0.011	0.875±0.023	0.944±0.009	0.968±0.003	0.995±0.002
	5	0.795±0.019	0.844±0.022	0.436±0.067	0.747±0.028	0.851±0.018	0.919±0.022	0.975±0.003	0.979±0.003	0.907±0.016	0.964±0.010	0.975±0.003	0.989±0.005
	6	0.798±0.041	0.750±0.125	0.348±0.039	0.435±0.035	0.914±0.023	0.790±0.018	0.969±0.008	0.960±0.023	0.879±0.003	0.918±0.006	0.987±0.003	0.967±0.011
	7	0.704±0.117	0.872±0.015	0.525±0.122	0.657±0.016	0.838±0.009	0.966±0.006	0.948±0.022	0.978±0.005	0.899±0.031	0.944±0.004	0.965±0.002	0.993±0.001
	8	0.856±0.018	0.728±0.021	0.602±0.087	0.356±0.027	0.734±0.001	0.711±0.038	0.979±0.003	0.953±0.004	0.939±0.017	0.879±0.009	0.952±0.000	0.963±0.006
	9	0.853±0.019	0.813±0.015	0.446±0.052	0.522±0.022	0.896±0.009	0.955±0.011	0.978±0.003	0.973±0.004	0.912±0.006	0.909±0.018	0.981±0.002	0.993±0.002
	MEAN	0.810±0.009	0.833±0.001	0.548±0.053	0.563±0.004	0.851±0.010	0.902±0.016	0.966±0.008	0.972±0.011	0.916±0.013	0.929±0.002	0.971±0.002	0.980±0.007
8-shot	0	0.855±0.038	0.916±0.013	0.631±0.039	0.521±0.014	0.833±0.015	0.978±0.011	0.979±0.006	0.988±0.001	0.932±0.005	0.931±0.002	0.960±0.003	0.995±0.001
	1	0.978±0.003	0.997±0.001	0.747±0.026	0.443±0.019	0.962±0.008	0.918±0.002	0.997±0.001	0.999±0.000	0.963±0.008	0.917±0.004	0.992±0.001	0.999±0.001
	2	0.661±0.010	0.803±0.055	0.651±0.063	0.584±0.032	0.790±0.030	0.870±0.014	0.944±0.002	0.963±0.013	0.947±0.006	0.929±0.002	0.956±0.006	0.965±0.002
	3	0.852±0.020	0.908±0.030	0.712±0.091	0.775±0.044	0.939±0.035	0.978±0.000	0.980±0.003	0.986±0.005	0.956±0.013	0.978±0.007	0.989±0.007	0.997±0.000
	4	0.808±0.012	0.885±0.039	0.353±0.038	0.608±0.025	0.856±0.017	0.978±0.003	0.971±0.002	0.982±0.006	0.930±0.019	0.953±0.005	0.971±0.003	0.997±0.000
	5	0.819±0.016	0.885±0.020	0.453±0.073	0.748±0.036	0.866±0.012	0.920±0.038	0.978±0.002	0.984±0.003	0.907±0.034	0.974±0.013	0.978±0.002	0.989±0.008
	6	0.834±0.005	0.850±0.019	0.352±0.025	0.469±0.039	0.925±0.008	0.855±0.036	0.975±0.001	0.977±0.003	0.881±0.002	0.925±0.005	0.988±0.001	0.990±0.006
	7	0.777±0.078	0.899±0.031	0.558±0.119	0.655±0.017	0.856±0.009	0.968±0.005	0.961±0.014	0.982±0.008	0.877±0.027	0.955±0.004	0.969±0.002	0.996±0.001
	8	0.870±0.020	0.767±0.023	0.555±0.085	0.354±0.023	0.744±0.005	0.777±0.011	0.981±0.003	0.958±0.006	0.880±0.016	0.885±0.008	0.953±0.001	0.968±0.002
	9	0.839±0.022	0.850±0.019	0.459±0.048	0.515±0.028	0.867±0.006	0.961±0.009	0.976±0.003	0.974±0.003	0.916±0.005	0.919±0.014	0.982±0.001	0.991±0.001
	MEAN	0.829±0.009	0.876±0.004	0.547±0.063	0.567±0.007	0.867±0.007	0.920±0.003	0.974±0.002	0.979±0.001	0.919±0.018	0.937±0.004	0.974±0.001	0.989±0.001

Table 11. Fine-grained AUROC and AUPRC results(mean±std) on One-vs-all protocol of MNIST datasets under various few-shot AD settings. Best results and the second-best results are respectively highlighted in red and blue.

Data subset		AUROC						AUPRC					
		SPADE	Patchcore	RegAD	CoOp	WinCLIP	Ours (InCTRL)	SPADE	Patchcore	RegAD	CoOp	WinCLIP	Ours (InCTRL)
2-shot	airplane	0.822±0.053	0.566±0.065	0.439±0.043	0.437±0.014	0.917±0.005	0.938±0.006	0.974±0.009	0.915±0.014	0.865±0.011	0.886±0.005	0.988±0.002	0.993±0.001
	automobile	0.929±0.013	0.499±0.041	0.681±0.026	0.800±0.007	0.932±0.002	0.945±0.005	0.991±0.002	0.898±0.012	0.946±0.004	0.969±0.003	0.992±0.000	0.995±0.001
	bird	0.721±0.093	0.583±0.040	0.413±0.019	0.424±0.003	0.925±0.002	0.931±0.011	0.952±0.019	0.923±0.010	0.880±0.016	0.886±0.003	0.990±0.002	0.992±0.001
	cat	0.748±0.087	0.567±0.010	0.516±0.027	0.359±0.003	0.942±0.000	0.940±0.002	0.956±0.019	0.918±0.002	0.914±0.012	0.875±0.002	0.992±0.000	0.992±0.001
	deer	0.869±0.037	0.738±0.034	0.657±0.021	0.436±0.006	0.928±0.005	0.943±0.002	0.978±0.008	0.956±0.006	0.933±0.012	0.884±0.006	0.990±0.001	0.993±0.000
	dog	0.716±0.071	0.529±0.034	0.588±0.024	0.533±0.028	0.842±0.003	0.856±0.009	0.947±0.014	0.911±0.017	0.933±0.003	0.914±0.002	0.978±0.000	0.981±0.001
	frog	0.832±0.026	0.791±0.018	0.650±0.013	0.288±0.045	0.920±0.004	0.931±0.002	0.974±0.004	0.967±0.004	0.946±0.005	0.867±0.008	0.990±0.000	0.992±0.000
	horse	0.787±0.102	0.621±0.027	0.483±0.015	0.762±0.026	0.916±0.001	0.936±0.007	0.962±0.021	0.930±0.007	0.896±0.015	0.957±0.006	0.987±0.000	0.990±0.001
	ship	0.894±0.020	0.601±0.028	0.519±0.027	0.549±0.005	0.951±0.001	0.956±0.006	0.986±0.002	0.929±0.008	0.898±0.004	0.926±0.002	0.992±0.000	0.996±0.001
	truck	0.906±0.058	0.529±0.071	0.391±0.054	0.679±0.013	0.978±0.002	0.976±0.004	0.988±0.008	0.912±0.008	0.878±0.009	0.947±0.005	0.997±0.001	0.998±0.001
	MEAN	0.823±0.014	0.602±0.009	0.534±0.005	0.527±0.011	0.925±0.001	0.935±0.002	0.971±0.003	0.926±0.002	0.909±0.003	0.911±0.002	0.990±0.001	0.992±0.000
4-shot	airplane	0.831±0.038	0.578±0.025	0.474±0.036	0.471±0.008	0.919±0.002	0.942±0.011	0.974±0.007	0.918±0.007	0.886±0.007	0.905±0.004	0.989±0.001	0.995±0.004
	automobile	0.929±0.018	0.538±0.025	0.728±0.018	0.808±0.005	0.933±0.000	0.948±0.022	0.991±0.002	0.909±0.006	0.951±0.002	0.971±0.003	0.992±0.000	0.995±0.007
	bird	0.754±0.048	0.621±0.023	0.424±0.020	0.437±0.005	0.927±0.001	0.933±0.005	0.959±0.009	0.933±0.005	0.874±0.012	0.892±0.001	0.991±0.000	0.992±0.001
	cat	0.760±0.075	0.584±0.021	0.529±0.022	0.356±0.002	0.942±0.001	0.945±0.012	0.959±0.015	0.923±0.005	0.912±0.009	0.872±0.002	0.992±0.000	0.992±0.003
	deer	0.905±0.020	0.789±0.013	0.597±0.013	0.455±0.005	0.932±0.003	0.956±0.003	0.984±0.004	0.965±0.002	0.920±0.006	0.891±0.002	0.990±0.001	0.994±0.000
	dog	0.665±0.066	0.526±0.050	0.554±0.017	0.527±0.024	0.841±0.002	0.866±0.019	0.938±0.015	0.911±0.017	0.914±0.005	0.911±0.003	0.978±0.000	0.984±0.007
	frog	0.862±0.013	0.828±0.011	0.513±0.004	0.296±0.039	0.923±0.001	0.935±0.009	0.979±0.003	0.973±0.002	0.920±0.003	0.871±0.005	0.991±0.000	0.993±0.005
	horse	0.846±0.040	0.694±0.013	0.524±0.012	0.763±0.018	0.916±0.001	0.939±0.012	0.974±0.009	0.946±0.003	0.912±0.010	0.957±0.005	0.987±0.000	0.987±0.002
	ship	0.903±0.027	0.663±0.034	0.496±0.006	0.557±0.007	0.953±0.001	0.958±0.016	0.988±0.003	0.942±0.007	0.904±0.003	0.930±0.001	0.993±0.000	0.991±0.006
	truck	0.901±0.049	0.567±0.053	0.498±0.015									