

Supplementary Materials of “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

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 [1], VisA [14], ELPV [4], SDD [12], AITEX [11]), two medical image datasets (BrainMRI [10], HeadCT [10]), and two semantic anomaly detection datasets: MNIST [8] and CIFAR-10 [7] under both one-vs-all and multi-class protocols [3].

To assess the Generalist 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 1 provides the data statistics of MVTec AD and VisA, while Table 2 shows the test set statistics of the rest datasets.

A.2. Industrial Defect Inspection Datasets

MVTec AD [1] 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 [14] 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

Dataset	Subset	Type	Original Training	Original Test	
			Normal	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 1. 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.

various flaws, including surface defects such as scratches, dents, color spots or crack, and structural defects like misplacement or missing parts.

ELPV [4] 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 [12] 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

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

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 2. 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.

study, we only use its test set for evaluation.

AITEX [11] 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 [10] 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 [10] 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 [8] 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 normal class label is set to ‘even_number’.

CIFAR-10 [7] 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 normal class label is set to ‘animal’.

B. Implementation Details

B.1. Data Pre-processing

By default, for all CLIP-based models, including WinCLIP [6], CoOp [13], and InCTRL, we adopt the same CLIP implementation, OpenCLIP [5], 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 → 896), text embedding (512 → 640) and input resolution (224×224 → 240×240). 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 [6] except

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 3. Examples of normal and abnormal text prompts used in InCTRL and WinCLIP. $[c]$ represents a class label.

CIFAR-10 [14]. Table 3 provides several normal and abnormal examples of text prompts used in InCTRL (see [6] for the full list of the text prompts). The WinCLIP prompts fail to work for natural images like CIFAR-10, so the normal 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 4.

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 4 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

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 4. Number of Parameters and Per-image Inference time.

	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 5. Comparison with Traditional Medical AD Methods.

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 images on two other medical datasets (Brain Tumor MRI¹ and LAG [9]) in Table 5. The results of traditional medical anomaly detection methods are from Cai *et al.* [2] 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 5, 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 6 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 7 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.

¹The dataset is available at <https://www.kaggle.com/datasets/masoudnickparvar/brain-tumor-mri-dataset>.

