

Research Article

# Vision-Language In-Context Learning Driven Few-Shot Visual Inspection Model

Shiryu Ueno<sup>1,2</sup>, Yoshikazu Hayashi<sup>1</sup>, Shunsuke Nakatsuka<sup>1,2</sup>, Yusei Yamada<sup>1</sup>, Hiroaki Aizawa<sup>2</sup>

1. Faculty of Engineering, Gifu University, Japan; 2. Graduate School of Advanced Science and Engineering, Hiroshima University, Japan

We propose general visual inspection model using Vision-Language Model (VLM) with few-shot images of non-defective or defective products, along with explanatory texts that serve as inspection criteria. Although existing VLM exhibit high performance across various tasks, they are not trained on specific tasks such as visual inspection. Thus, we construct a dataset consisting of diverse images of non-defective and defective products collected from the web, along with unified formatted output text, and fine-tune VLM. For new products, our method employs In-Context Learning, which allows the model to perform inspections with an example of non-defective or defective image and the corresponding explanatory texts with visual prompts. This approach eliminates the need to collect a large number of training samples and re-train the model for each product. The experimental results show that our method achieves high performance, with MCC of 0.804 and F1-score of 0.950 on MVTec AD in a one-shot manner. Our code is available at <https://github.com/ia-gu/Vision-Language-In-Context-Learning-Driven-Few-Shot-Visual-Inspection-Model>.

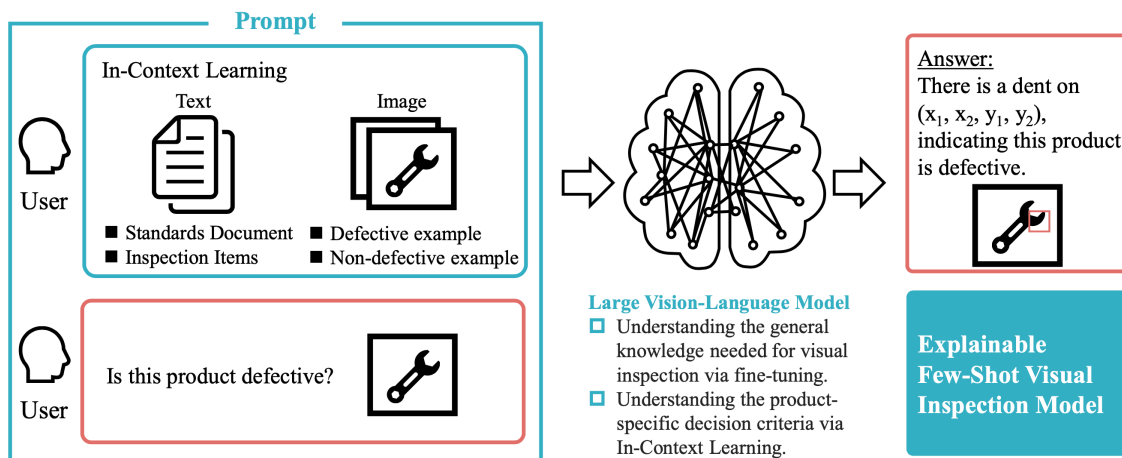
## 1. Introduction

In this study, we propose a method that can detect defective locations in new product images by using Vision-Language Model (VLM) <sup>[1][2]</sup> and In-Context Learning (ICL) <sup>[3][4]</sup>.

With the advancements in deep learning technology, the automation of visual inspection has become increasingly common in recent years. However, current visual inspection models inspect specific products by collecting a large number of images of the target product and training the model. Thus, these models are only applicable to the target products on which they have been trained, and re-training is necessary for new products. Although some methods can inspect multiple products with a

single model, they still require hyperparameter tuning or additional training for each product. In this study, we propose a general visual inspection model that leverages VLM and ICL allowing the inspection of new products without any hyperparameter tuning or model training.

Many of the current VLMs [5][6] leverage Large Language Model (LLM) to align visual and language features, demonstrating excellent performance in a wide range of tasks. These tasks range from basic image recognition tasks, such as classification, to advanced vision-language tasks, such as Visual Question Answering (VQA). However, these VLMs are not trained on specific tasks such as visual inspection.



**Figure 1.** Framework of our proposed method. We utilize ICL for multiple image inputs to give VLM the inspection criteria of new products. Our framework gives the coordinates of the defective location, which helps the user understand the model’s decision. In addition, it is easy to address by replacing the foundational model when a better VLM is proposed.

In this study, we propose a general visual inspection model that can detect defective locations in new products without any hyperparameter tuning or model re-training, using VLM and ICL. The framework of our proposed method is shown in Fig. 1. First, we fine-tune the VLM for general visual inspection with a dataset constructed from a diverse set of non-defective and defective product images collected from the web. In this study, we use ViP-LLaVA [7], which has been trained on visual prompt recognition, as the foundation of our VLM, and fine-tune it with our dataset. In addition, in typical visual inspection processes by humans, inspectors use inspection standards for the target products. To emulate this inspection process by human, we use ICL during the evaluation to provide

an example of non-defective or defective product image along with explanatory texts that serve as inspection criteria. ICL is a method the model learns from few-shot input-output examples as prompts, without parameter update. Using ICL during the inference of new products, we provide VLM with inspection criteria, enabling specific inspection of target products. Since ICL performance varies significantly based on the provided examples, we propose an algorithm that can select high-quality example based on the distance in Euclidean space. Consequently, our proposed method does not need to collect a large number of images or to re-train the model for each target product.

In summary, our main contributions are:

- We propose a general visual inspection model capable of inspections and detecting defective locations for new products using VLM and ICL with only an example. In our proposed method, fine-tune VLM on visual inspection and utilize ICL enabling the inspection of specific products.
- We construct a new dataset consisting of diverse non-defective and defective products collected from the web, along with unified formatted output, for fine-tuning. Also, our dataset includes coordinates of defective locations for defective products, ensuring the explainability of the model.
- To empirically verify the proposed methodology, we evaluate on MVTec AD<sup>[8]</sup> and VisA<sup>[9]</sup>. Our method achieves MCC<sup>[10]</sup> of 0.804 and F1-score<sup>[11]</sup> of 0.950 on the MVTec AD dataset in a one-shot manner.

## 2. Related Work

### 2.1. Visual Inspection

Many visual inspection methods based on deep learning are trained only on non-defective images<sup>[12]</sup> [13]. Thus, such methods require the collection of training samples and the re-training of the model for each target product. Consequently, it is challenging to apply the same model to different products without re-training.

Recently, visual inspection methods combining vision and language have been proposed. AnomalyGPT<sup>[14]</sup> can detect defective locations by learning an image decoder from non-defective and pseudo-defective images. However, AnomalyGPT utilizes PaDiM or PatchCore<sup>[15]</sup> for anomaly maps, and these methods need re-training for each products. WinCLIP<sup>[16]</sup> calculates the similarity between images and texts of non-defective and defective images using CLIP<sup>[17]</sup> and can detect defective

locations by using relative anomaly scores. However, WinCLIP only assigns anomaly scores to test samples during inference. To inspect correctly, it is necessary to experimentally determine the optimal threshold on test samples. Thus, these existing approaches cannot be considered general visual inspection models.

## 2.2. Vision-Language Model

VLMs leverage LLM to align visual and language features, demonstrating excellent performance across a wide range of tasks, from basic image recognition tasks such as classification, to advanced vision-language tasks, such as VQA. For example, LLaVA<sup>[18]</sup> inputs the vision embedding vectors and language embedding vectors into the LLM decoder to learn the alignment between vision and language. LLaVA has spawned many derivative methods, among which ViP-LLaVA focuses on visual prompt recognition by utilizing a dataset where arrows or visual cues are directly embedded in the input images, thereby strengthening the alignment between low-level image details and language. However, these VLMs have not been trained on visual inspection tasks and thus lack the general knowledge for visual inspection<sup>[2]</sup>.

## 2.3. In-Context Learning

ICL is a method that the model learns from few-shot input-output examples as prompts, without updating model parameters. For instance, given the input “Example input: (4, 2), Example output: 6, Question: (5, 6),” the model infers from the provided example that the task is addition and can answer “11.” In multi-modal ICL, the model makes inferences based on images, prompts, and their examples. Many VLMs are trained on diverse image-text pairs, enabling them to acquire ICL capabilities<sup>[19]</sup>.

Some VLMs are explicitly built to enhance ICL capabilities. Otter<sup>[20]</sup> enhances ICL capabilities by fine-tuning Open Flamingo<sup>[21]</sup> on MIMIC-IT<sup>[22]</sup>, which is in an ICL and Instruction Tuning format. At the same time, not to forget the knowledge of Open Flamingo, Otter only update parameters of Perceiver Resampler and Cross Attention Layer in language model. Similarly, LCL<sup>[23]</sup> proposes a new evaluation dataset, ISEKAI, which includes new concepts in the examples, making it challenging without seeing the examples. To address ISEKAI, LCL enhances its ICL capability by fully fine-tuning Shikra<sup>[6]</sup> on a custom dataset based on ImageNet<sup>[24]</sup>. However, in practice, these VLM explicitly designed to enhance ICL capabilities do not necessarily outperform regular VLM<sup>[19],[4]</sup>.

## 3. Proposed Method

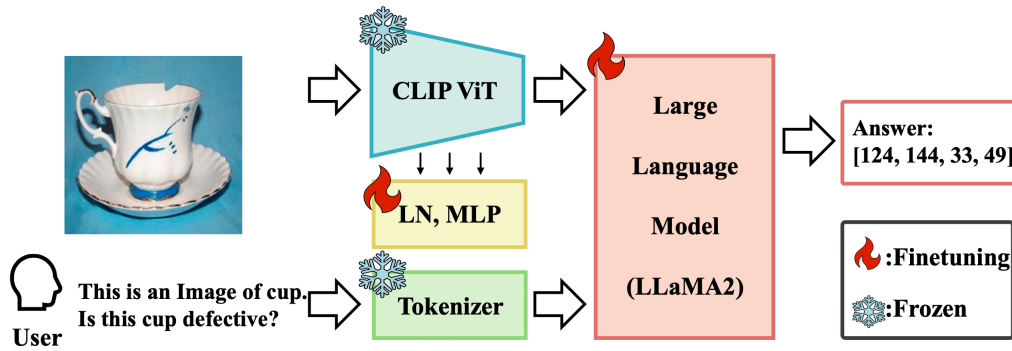
### 3.1. Overview

In this study, we propose a general visual inspection model that combines VLM and ICL, enabling the specific inspection of new products without parameter optimization. In addition, by constructing unified output format dataset for fine-tuning, we enable quantitative evaluation of visual inspections using VLM. An overview of the proposed method is shown in Fig. 1.

### 3.2. Model

In this study, we use ViP-LLaVA<sup>[7]</sup> as the foundational VLM. ViP-LLaVA is a model that improves recognition capabilities for visual prompts by fine-tuning LLaVA 1.5<sup>[5]</sup> on a dataset where red circles or arrows are overlaid on the original images. In addition to this, ViP-LLaVA utilizes the multi-level visual features to address the tendency of CLIP's deeper features to overlook low-level details. These methodologies improves the recognition capability for low-level details, which is especially needed for visual inspection. ViP-LLaVA has not been trained on visual inspection tasks.

The model architecture of ViP-LLaVA is shown in Fig. 2. ViP-LLaVA consists of a vision encoder to extract visual features, LayerNorm<sup>[25]</sup> and an MLP to tokenize visual features, a tokenizer to tokenize the language, and an LLM to generate text from these tokens. The vision encoder is CLIP-ViT-L/14<sup>[17]</sup>, and the LLM is LLaMA2<sup>[26]</sup> During fine-tuning, we update the parameters of the LayerNorm, MLP, and LLM in accordance with the ViP-LLaVA procedure.



**Figure 2.** Architecture of ViP-LLaVA. After providing an image and the corresponding text, the image is tokenized by CLIP ViT, LayerNorm, and MLP layers, while the text is tokenized by tokenizer. Then the visual tokens and the text tokens are given to the LLM to generate the answer.

### 3.3. Dataset for Fine-tune

To enhance the general knowledge of existing VLM for visual inspection, we collect images of non-defective and defective products from the web. The image collection process consists of five main steps:

1. Generate product names and inspection-related keywords (e.g., “disk”, “broken disk”, “discolored disk”) by GPT-4<sup>[27]</sup>.
2. Expand the keywords into eight languages: English, Chinese, Spanish, French, Portuguese, German, Italian, and Japanese.
3. Perform image searches using the expanded keywords and collect images with selenium.
4. Remove duplicate or unclear images.
5. Annotate the defective location coordinates for the defective images in the remaining set.

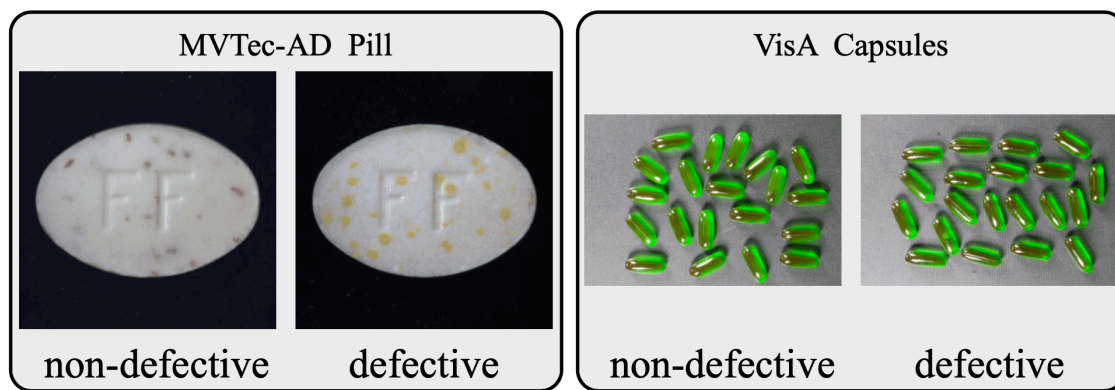
Through this procedure, we collect images of various products. Each product category includes images of non-defective and defective products with up to five types of defects. Finally, we obtained a final set of 941 images of 84 categories.

After collecting the images, we construct a dataset for fine-tuning. The format of the dataset is based on VQA (i.e., a pair of question and answer for each images). Question is “This is an image of {product}. Does this {product} in the image have any defects? If yes, please provide the anomaly mode

and the bounding box coordinate of the region where the defect is located. If no, please say None.”, answer is coordinates of the defective location for defective image, “None” for non-defective image.

### 3.4. In-Context Learning Driven Visual Inspection

It is challenging to inspect new products from a single image. An example is shown in Fig. 3. As shown in Fig. 3, some products need their specific inspection criteria for accurate visual inspection. Thus, in this study, we utilize ICL to provide an example of non-defective or defective image along with explanatory texts that serve as inspection criteria. Based on the example, model precisely inspects new products.



**Figure 3.** Examples of non-defective and defective images of “Pill” in MVTEC AD and “Capsules” in VisA. For “Pill”, the non-defective image also contains red spots, making it difficult to inspect. Similarly, for “Capsules”, the non-defective image also contains brown stains.

In addition, in multi-modal ICL, example images significantly influence the output of VLM<sup>[28]</sup>. RICES<sup>[29]</sup> is an existing algorithm for selecting examples in ICL, it uses cosine similarity of features. However, cosine similarity can yield high values when the scales of features differ or when the feature dimensions are large, failing to accurately evaluate similarity<sup>[30]</sup>. Thus, in this study, we propose a new selection algorithm. Our proposed method is shown in Eq. (1).

$$\operatorname{argmin} (\min \|f(x_i) - f(x_q)\|_2) \quad (1)$$

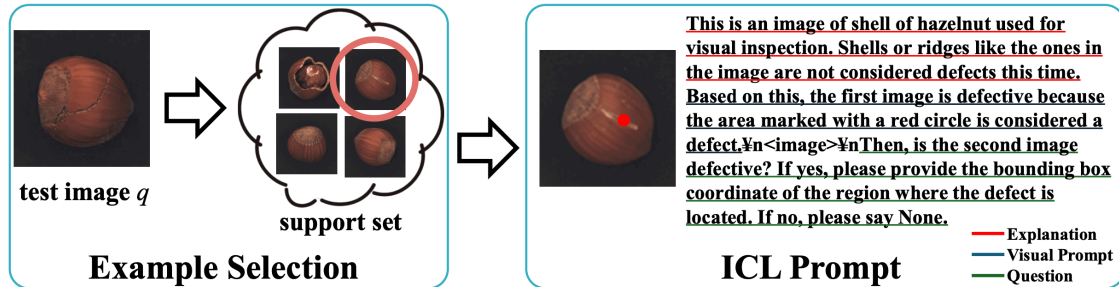
Where  $x$  denotes the image,  $f$  denotes the vision encoder (pre-trained ResNet50<sup>[31]</sup>), and  $q$  denotes the index of the test image for inference,  $i$  denotes the index of the image except for  $q$ . Eq. (1) is an algorithm that selects neighboring image of the test image as example based on Euclidean distance. By

this, Eq. (1) takes into account the scale and dimensions of the features, and expected to select better example compared to RICES.

## 4. Experiment

### 4.1. Settings

The dataset used for fine-tuning was collected and created as described in Sec. 3.3. We perform one-stage fine-tuning for 300 epochs using ZeRO2 with XTuner<sup>[32]</sup> on 8 NVIDIA 6000Ada-48GB GPUs. The batch size is set to 4, thus the global batch size is set to 32. We utilized the AdamW<sup>[33]</sup> optimizer and  $1e-4$  learning rate, with the warm-up ratio set to 0.03. We also apply cosine decay to the learning rate. To evaluate the performance, we used MVTEC AD and VisA. These datasets were not used at all during training. During the evaluation, one example image is given by using ICL (i.e., one-shot manner). The method for selecting the example image follows the procedure described in Eq. (1). Also, to evaluate the effectiveness of the proposed example image selection algorithm, we compare its performance without example image and with that of RICES. Finally, framework of the evaluation is shown in Fig. 4.



**Figure 4.** Framework of evaluation. First, select the example based on Eq. (1), then infer the test image with ICL.

We use F1-score<sup>[11]</sup> and Matthews Correlation Coefficient (MCC)<sup>[10]</sup> for the evaluation. F1-score, as shown in Eq. (2), is a common metric for binary classification.

$$\text{F1-score} = \frac{2 \times \text{TP}}{2 \times \text{TP} + \text{FP} + \text{FN}} \quad (2)$$

As shown in Eq. (2), F1-score does not use prediction of true negative. Thus, when there is a large number of positive samples during inference, the performance can be significantly inflated by



predicting all samples as positive. For instance, in MVTec AD, with 1,258 positive samples and 4,67 negative samples in the test data, F1-score shows a high value of 0.844. This shows F1-score is suspicious when there is a bias in the test data. Thus, we use not only F1-score but also MCC as evaluation metrics. MCC is reported to be adequate for binary classification, particularly for better consistency and less variance<sup>[34][35]</sup>. MCC is shown in Eq. (3).

$$\text{MCC} = \frac{\text{TP} \times \text{TN} - \text{FP} \times \text{FN}}{\sqrt{(\text{TP} + \text{FP})(\text{TP} + \text{FN})(\text{TN} + \text{FP})(\text{TN} + \text{FN})}} \quad (3)$$

MCC ranges from -1 to 1, where 1 indicates perfect prediction of all samples, -1 indicates incorrect prediction of all samples, and 0 indicates random prediction. In the previously mentioned example, MCC cannot be calculated because the denominator becomes zero. Thus, in this study, we use both F1-score and MCC for evaluation. The evaluation methods include assessing the performance for each product within each dataset, as well as the overall performance across the entire dataset.

## 4.2. Evaluation of Results

### 4.2.1. Result of MVTec AD

The results for MVTec AD are shown in Table 1. The settings are as follows: “Vanilla” for ViP-LLaVA before fine-tune, “w/o ICL” for ViP-LLaVA after fine-tune without using an example during inference, “ICL (RICES)” for using a selected example image with the RICES algorithm during inference, and “ICL (Ours)” for using a selected example image with Eq. (1) during inference. In each settings, results of F1-score and MCC are in a row. From the table, we confirm that providing an example significantly improves performance. This demonstrates the effectiveness of our framework. Additionally, compared to RICES, our selection algorithm achieves improvement in performance with an increase in MCC, demonstrating the effectiveness of our algorithm.

Settings	Vanilla		w/o ICL		ICL (RICES)		ICL (Ours)	
	F1-score	MCC	F1-score	MCC	F1-score	MCC	F1-score	MCC
Bottle	N/A	N/A	0.863	N/A	0.892	0.510	0.917	<b>0.610</b>
Cable	N/A	N/A	0.400	0.338	0.795	0.564	0.899	<b>0.754</b>
Capsule	N/A	N/A	0.750	0.426	0.912	0.384	0.946	<b>0.658</b>
Carpet	0.044	0.074	1.000	<b>1.000</b>	0.983	0.929	1.000	<b>1.000</b>
Grid	N/A	N/A	0.973	0.910	0.884	0.476	0.982	<b>0.935</b>
Hazelnut	0.228	0.226	0.780	N/A	0.795	0.257	0.800	<b>0.289</b>
Leather	0.043	0.076	1.000	<b>1.000</b>	1.000	<b>1.000</b>	1.000	<b>1.000</b>
Metal Nut	N/A	N/A	0.832	0.540	0.912	0.468	0.989	<b>0.947</b>
Pill	N/A	N/A	0.838	0.402	0.922	0.368	0.968	<b>0.814</b>
Screw	0.209	0.140	0.851	0.244	0.903	0.506	0.925	<b>0.673</b>
Tile	N/A	N/A	0.957	0.870	0.977	<b>0.916</b>	0.977	<b>0.916</b>
Toothbrush	N/A	N/A	0.906	0.633	0.866	0.418	0.921	<b>0.697</b>
Transistor	N/A	N/A	0.762	0.592	0.871	0.780	0.894	<b>0.821</b>
Wood	N/A	N/A	0.976	0.752	0.992	<b>0.965</b>	0.992	<b>0.965</b>
Zipper	N/A	N/A	0.741	0.541	0.975	0.879	0.987	<b>0.941</b>
All category	0.042	0.068	0.860	0.519	0.917	0.665	0.950	<b>0.804</b>

**Table 1.** Result of MVTEC-AD. 'N/A' means that zero division occurred. Bold means the highest performance.

Next, for qualitative evaluation, the visualization of the model prediction is shown in Fig. 5. As shown in the figure, our approach can roughly detect defective locations, which means the model recognizes the defects in the image. However, the model cannot detect multiple defects or logical defects, such as

those in “Cable”. This is due to the lack of variety in the training dataset. Thus, further image collection and an enlarged training dataset are required for performance improvement.

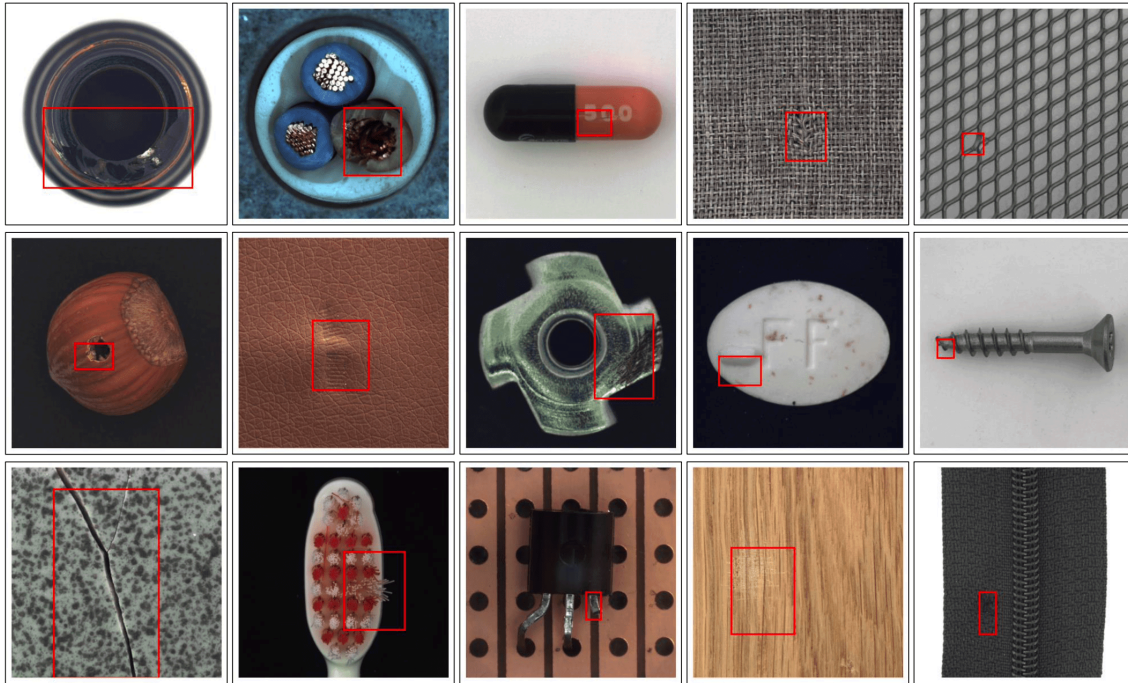


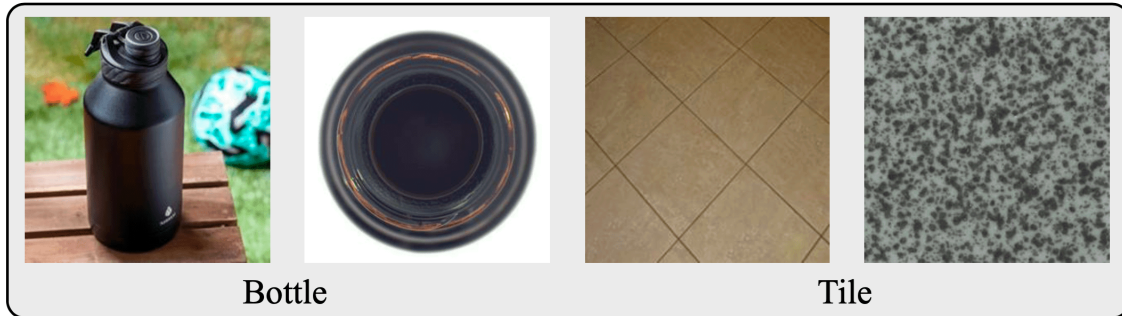
Figure 5. Visualize the model prediction for MVTec AD.

Also, for some products like “Hazelnut”, while our approach improved the performance, it is still insufficient for real-world conditions. For “Hazelnut”, the model detected thin parts as defective, indicating that the model does not fully leverage ICL. Thus, providing detailed inspection criteria is necessary. It has been reported that increasing the number of examples improves ICL performance<sup>[36]</sup><sup>[37]</sup>. Alternatively, further performance improvement is expected by proposing an optimal selection algorithm that selects multiple example images (based on the query strategies, including those from Deep Active Learning<sup>[38]</sup><sup>[39]</sup>).

Additionally, for all products, although coordinates are output, their positions deviate from the actual defective locations. Indeed, pixel-level AUROC was 0.730, which is very low compared to the existing methods. This is because the CrossEntropyLoss used for training uniformly calculates the loss for differences in token values. For example, when the ground truth of the starting x-coordinate is 100, the loss is the same when the model outputs 101 and 900 (assuming the prediction probabilities are equal). Thus, CrossEntropyLoss is not optimal for tasks requiring specific numerical outputs like

coordinates. However, existing VLMs are trained with CrossEntropyLoss, meaning their outputs are text-based cannot be safely converted to floats (with gradient flow intact), thus performance improvement is expected by constructing a multi-head VLM for defect detection and modifying the loss function to alternatives like Mean Squared Error or GIoU Loss.

While “Bottle” has the same product in the training dataset, their performance is lower compared to “Wood”, which also has the same product in the training dataset. This is likely due to the significant differences in appearance between the images in the training data and those in MVTec AD, as shown in Fig. 6. However, despite the differences in appearance, “Tile” shows high performance, confirming the generalization capability for some products. Also, to prevent the forgetting of knowledge acquired during pre-training when fine-tuning, it is necessary to use Parameter Efficient Fine-Tuning methods, such as Low-Rank Adaptation<sup>[40]</sup>, which forget less than fine-tuning<sup>[41]</sup>.



**Figure 6.** Examples of the images of “Bottle”, and “Tile” from the collected images and MVTec AD.

#### 4.2.2. Result of VisA

The results for VisA are shown in Table 2. The table follows the same format as Table 1. From the table, it can be confirmed that the performance improves by using ICL in VisA as well, demonstrating the effectiveness of the proposed framework. However, compared to RICES, our selection algorithm does not show significant improvement. This is because both RICES and our selection algorithm are based on similarity, which depends on the data distribution. Most of the products in VisA are too widely distributed (e.g., “Macaroni”, “PCB”). Thus, proposing a more distribution-robust selection algorithm could potentially improve performance. Also, it can be seen that the performance does not improve regardless of the presence of ICL when there are two or more products in the image, especially if those products are not aligned. In fact, “Macaroni1”, which is neatly aligned, shows

higher qualitative and quantitative performance compared to “Macaroni2”, which is randomly arranged. This is likely due to the lack of training dataset that considers differences in product positions and orientations. Thus, performance improvement is expected by collecting fine-tuning data and performing data augmentation, such as rotation and flipping. Simultaneously, it should be noted that for some products, positional shifts or orientation differences may be defined as defects.

Settings	Vanilla		w/o ICL		ICL (RICES)		ICL (Ours)	
Product Name	F1-score	MCC	F1-score	MCC	F1-score	MCC	F1-score	MCC
Candle	N/A	N/A	0.635	<b>0.539</b>	0.692	0.241	0.694	0.253
Capsules	N/A	N/A	0.599	0.415	0.841	<b>0.513</b>	0.809	0.389
Cashew	N/A	N/A	0.814	0.623	0.890	0.670	0.889	<b>0.674</b>
Chewinggum	N/A	N/A	0.921	0.758	0.921	0.758	0.935	<b>0.804</b>
Fryum	N/A	N/A	0.867	0.699	0.917	<b>0.741</b>	0.888	0.648
Macaroni1	N/A	N/A	0.760	<b>0.502</b>	0.685	0.204	0.683	0.190
Macaroni2	N/A	N/A	0.669	<b>0.071</b>	0.667	N/A	0.667	N/A
PCB1	N/A	N/A	0.131	0.190	0.891	<b>0.792</b>	0.875	0.762
PCB2	N/A	N/A	0.347	0.343	0.772	<b>0.493</b>	0.763	0.471
PCB3	N/A	N/A	0.243	0.248	0.747	0.503	0.751	<b>0.513</b>
PCB4	N/A	N/A	0.622	0.516	0.801	0.594	0.817	<b>0.610</b>
Pipe Fryum	N/A	N/A	0.870	0.726	0.920	0.744	0.929	<b>0.774</b>
All category	N/A	N/A	0.671	0.429	0.800	<b>0.492</b>	0.795	0.479

Table 2. Result of VisA.

For qualitative evaluation, the visualization of the model prediction is shown in Fig. 7. As shown in Fig. 7, for products that have multiple objects like “Candle” or “Capsules”, the model prediction gets worse. As mentioned, our dataset is still insufficient for generalization because there are limited

products and they are mostly single object. In addition, images with multiple objects are highly distributed compared to the images with single object, which influences the performance of ICL because the selection algorithms depend on the distribution.

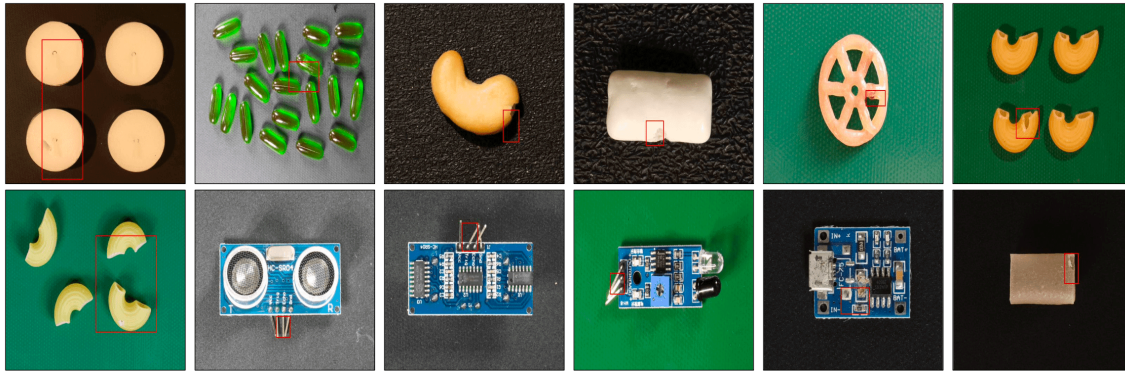


Figure 7. Visualize the model prediction for VisA.

## 5. Conclusion

In this study, we propose a general visual inspection model based on a few images of non-defective or defective products along with explanatory texts serving as inspection criteria. For future work, further performance improvement is expected by collecting more images for fine-tuning. In this study, we enabled visual inspection using VLM by training on a dataset consisting of only 941 images, which is very small compared to the pre-training dataset of VLM. Another consideration is to construct the multi-head VLM and change of the loss function. Furthermore, introducing the example image selection algorithm is another way for improvement. Specifically, existing algorithms are for selecting one example image for the inspection, so proposing an optimal selection algorithm for many example images improves model performance. In addition, the proposed method is based on VLM, so by adding the rationale statements for the decision in the response, model explainability is expected to improve, and performance could be enhanced through multitasking. Finally, in our study, we evaluate only on MVTec AD and VisA. However, more comprehensive benchmark such as MMAD<sup>[42]</sup> is required.

# Appendix

## 6. Product Category

As mentioned in Sec. 3.3, we fine-tuned LVLMM using a diverse set of non-defective and defective images of various products collected from the web to enhance the visual inspection capabilities of LVLMM. The product names in the dataset we used are shown in Fig. 8. Note that we renamed the products during dataset construction after collection (e.g., CD → disk, carpet → textiles). All the images used for training will be publicly available at <https://github.com/ia-gu/Vision-Language-In-Context-Learning-Driven-Few-Shot-Visual-Inspection-Model>.

- apple
- asphalt
- automobile
- bag
- battery
- battery holder
- biscuit
- bottle
- bowl
- cable
- cross section
- camera
- can
- chair
- cherry
- circuit board
- coin
- concrete
- corn
- CPU fan
- cup
- disk
- display
- egg
- eggplant
- fan
- fasteners
- fork
- gauge
- glass cup
- glasses
- grape
- hat
- keyboard
- knife
- LAN cable
- leaf
- LED
- lemon
- lens
- locker
- mango
- marble
- metal
- monitor
- motor
- nut
- orange
- pan
- pants
- paper
- peach
- pen
- ping pong ball
- plastic bag
- plate
- RAM
- road
- saddle
- shirt
- shoe sole
- shoes
- smartphone
- smartphone lens
- spoon
- spray
- SSD
- strawberry
- table
- tablet
- test tube
- textiles
- tile
- tomato
- toothbrush
- transistor
- valve
- videotape
- vinyl record
- vinyl siding
- wood
- wooden wheel
- zipper

Figure 8. Product category of our dataset.

## 7. Result of ViP-LLaVA Before Fine-tuning

An example of the prediction results on MVTec AD of ViP-LLaVA before fine-tuning is shown in Fig. 9. As illustrated, the vanilla ViP-LLaVA fails correct inspection, predicts both non-defective and defective products as defective. Moreover, the format of the response text is inconsistent, making it



challenging to perform a consistent quantitative evaluation. These results confirm the effectiveness of fine-tuning using our dataset. On the other hand, from the third result, it can be seen that the vanilla ViP-LLaVA is capable of describing the type and location of the defect with statements such as 'The defect is located at the center of the leather, and it appears to be a hole.' Therefore, by adding rationale statements to our dataset, it is suggested that the proposed method could output not only defective location coordinates but also the rationale statements for the judgment.





	<p><u>Question:</u> This is an image of bottle. Does this bottle in the image have any defects? If yes, please provide the bounding box coordinate of the region where the defect is located.</p> <p><u>ViP-LLaVA Prediction:</u> Yes, the bottle in the image has a defect. The defect is located at the top of the bottle, where the lid is missing.</p>
	<p><u>Question:</u> This is an image of bottle. Does this bottle in the image have any defects? If yes, please provide the bounding box coordinate of the region where the defect is located.</p> <p><u>ViP-LLaVA Prediction:</u> Yes, the bottle in the image has a defect. The defect is located at the top of the bottle, near the neck.</p>
	<p><u>Question:</u> This is an image of leather. Does this leather in the image have any defects? If yes, please provide the bounding box coordinate of the region where the defect is located.</p> <p><u>ViP-LLaVA Prediction:</u> Yes, the leather in the image has a defect. The defect is located at the center of the leather, and it appears to be a hole.</p>
	<p><u>Question:</u> This is an image of leather. Does this leather in the image have any defects? If yes, please provide the bounding box coordinate of the region where the defect is located.</p> <p><u>ViP-LLaVA Prediction:</u> The leather in the image has a defect. The defect is located in the upper left corner of the image.</p>

Figure 9. Result of the ViP-LLaVA before fine-tuning.

## 8. Ablation Study of In-Context Learning for MVTec AD

In our main experiments, we confirmed that selectively providing a single example image during ICL in the evaluation improves performance. Here, we use MVTec AD to verify the effectiveness of the proposed method by comparing results when example images are selected randomly. At the same time, we compare results when the number of example images is increased.



The experimental results are shown in Table 3. From the table, it can be seen that the highest performance is achieved when a single non-defective example is provided randomly (note that this is lower than the result of 'w/o ICL'). Additionally, simply increasing the number of provided examples does not improve performance; on the contrary, it decreases. This indicates that in ICL, the influence of the examples is significant, and increasing the number of examples without considering their relevance to the query image leads to performance degradation. Although our proposed method and RICES are algorithms specialized in selecting a single example image, the performance improvement expected from increasing the number of examples in ICL suggests that further improvements could be achieved by proposing an algorithm for selecting two or more examples.

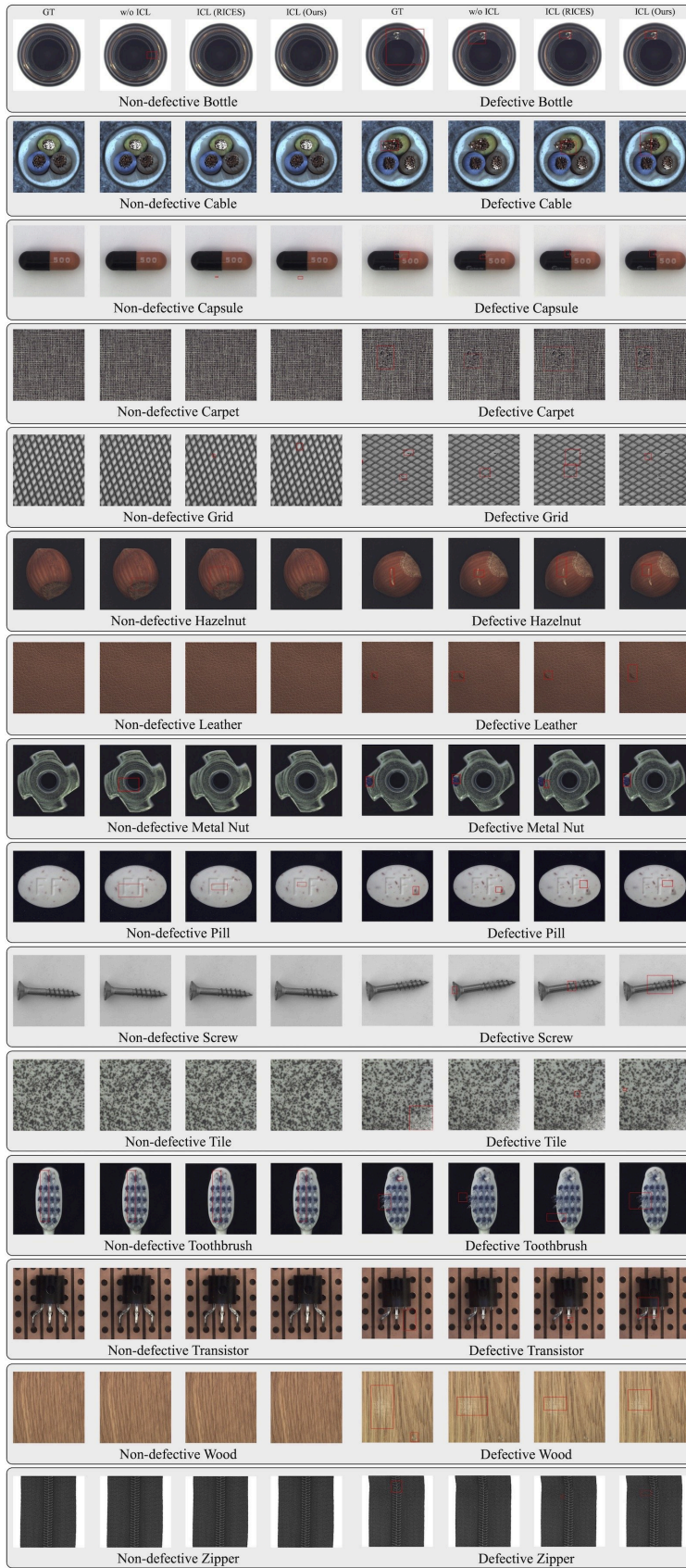
Settings	1-pos	1-neg	2-pos-pos	2-neg-neg	2-pos-neg	2-neg-pos
Bottle	0.509	0.589	N/A	N/A	0.201	0.498
Cable	0.124	0.589	0.184	0.305	0.275	0.056
Capsule	0.185	0.353	0.006	N/A	0.305	0.164
Carpet	0.195	0.892	0.356	0.931	0.678	0.124
Grid	0.598	0.816	0.547	0.838	0.684	0.544
Hazelnut	0.428	0.198	0.519	N/A	0.105	0.495
Leather	0.578	0.939	0.343	1.000	0.734	0.603
Metal Nut	0.226	0.402	0.002	0.427	0.237	0.244
Pill	0.374	0.440	0.227	0.282	0.405	0.313
Screw	0.123	0.084	0.155	0.004	0.167	0.156
Tile	0.724	0.764	0.621	0.790	0.790	0.689
Toothbrush	N/A	0.484	N/A	0.496	N/A	N/A
Transistor	0.205	0.556	0.004	0.510	0.275	0.151
Wood	0.421	0.896	0.280	0.849	0.775	0.421
Zipper	0.260	0.482	0.280	0.418	0.410	0.242
All category	0.279	0.475	0.197	0.455	0.360	0.267

**Table 3.** Result of the comparison of In-Context Learning. All the results are by MCC. Each of the index means “1-pos” gives one defective example image, “1-neg” gives one non-defective example image, “2-pos-pos” gives two defective example images, “2-neg-neg” gives two non-defective example images, “2-pos-neg” gives one defective example image and one non-defective example image in a row, “2-neg-pos” is vice versa.

## 9. Visualization Results of All Products for MVTec AD and VisA

The visualization of the model prediction for MVTec AD and VisA is shown in Fig. 10 and Fig. 11. We can see the same tendency that when there are two or more products in the image, the performance

decreases, and ICL does not work well. The performance decrease is due to the lack of dataset diversity. Most images in our dataset contain a single product. The reason why ICL does not work well is that when there are two or more products in each image of one category, the diversity within the category increases, and the algorithm for calculating similarity or distance fails to perform effectively. Thus, for future work, it is noted that simply calculating similarity or distance may fail in specific domains.



**Figure 10.** Visualize the model prediction of all products for MVTec AD.

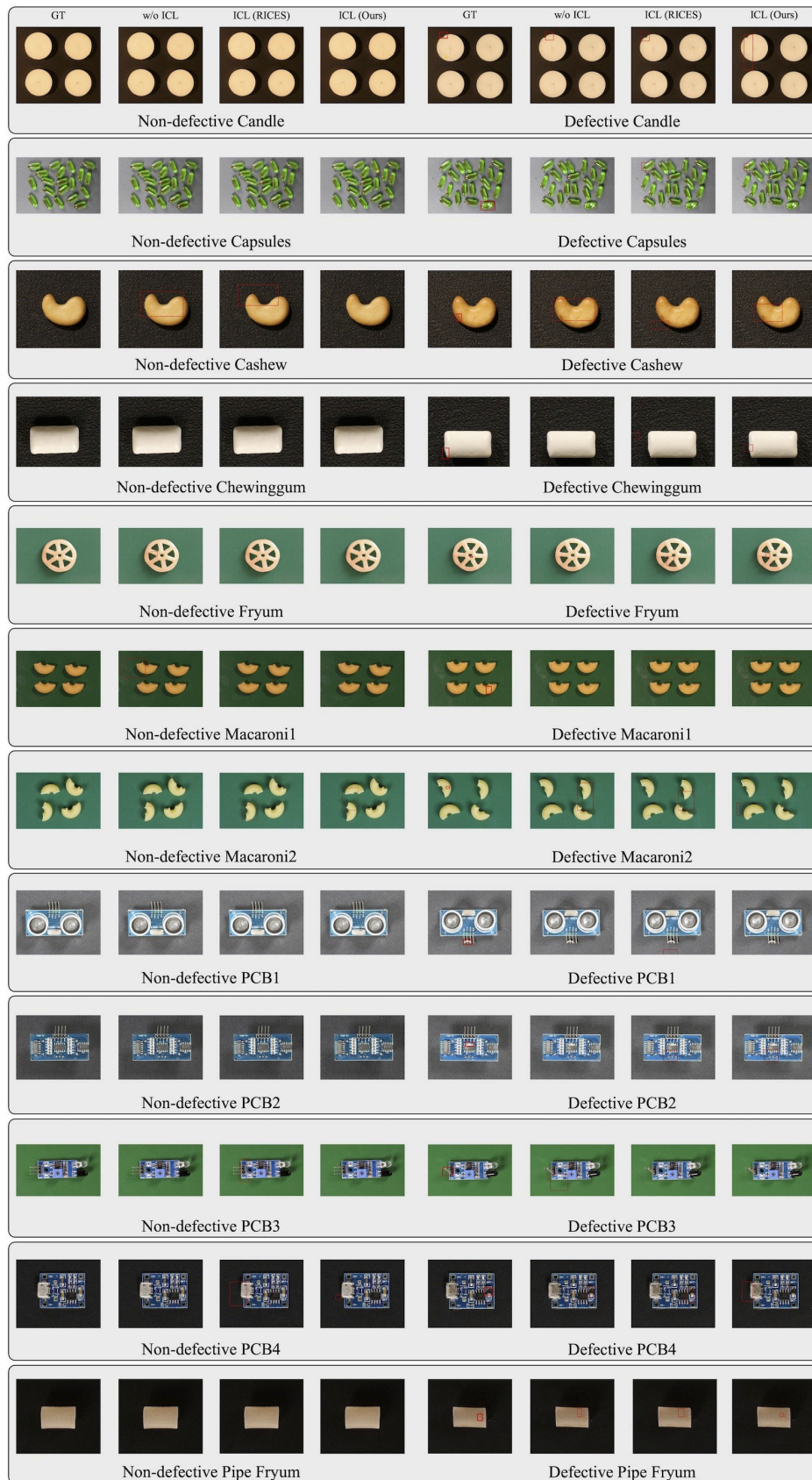


Figure 11. Visualize the model prediction of all products for VisA.

## References

1. <sup>^</sup>Yin S, Fu C, Zhao S, Li K, Sun X, Xu T, Chen E (2024). "A Survey on Multimodal Large Language Model s". arXiv. arXiv:2306.13549. Available from: <https://arxiv.org/abs/2306.13549>.
2. <sup>^</sup>Liu Y, Duan H, Zhang Y, Li B, Zhang S, Zhao W, Yuan Y, Wang J, He C, Liu Z, Chen K, Lin D (2024). "MBench: Is Your Multi-modal Model an All-around Player?" arXiv. arXiv:2307.06281. Available from: <https://arxiv.org/abs/2307.06281>.
3. <sup>^</sup>Dong Q, Li L, Dai D, Zheng C, Wu Z, Chang B, Sun X, Xu J, Li L, Sui Z (2023). "A Survey on In-context Learning". arXiv. arXiv:2301.00234. Available from: <https://arxiv.org/abs/2301.00234>.
4. <sup>^</sup>Zong Y, Bohdal O, Hospedales T (2024). "VL-ICL Bench: The Devil in the Details of Benchmarking Multimodal In-Context Learning". arXiv. arXiv:2403.13164. Available from: <https://arxiv.org/abs/2403.13164>.
5. <sup>^</sup>Liu H, Li C, Li Y, Lee YJ (2024). "Improved Baselines with Visual Instruction Tuning". arXiv. arXiv:2310.03744. Available from: <https://arxiv.org/abs/2310.03744>.
6. <sup>^</sup>Chen K, Zhang Z, Zeng W, Zhang R, Zhu F, Zhao R (2023). "Shikra: Unleashing Multimodal LLM's Referential Dialogue Magic". arXiv. arXiv:2306.15195. Available from: <https://arxiv.org/abs/2306.15195>.
7. <sup>^</sup>Cai M, Liu H, Park D, Mustikovela SK, Meyer GP, Chai Y, Lee YJ (2024). "ViP-LLaVA: Making Large Multimodal Models Understand Arbitrary Visual Prompts". arXiv. arXiv:2312.00784. Available from: <https://arxiv.org/abs/2312.00784>.
8. <sup>^</sup>Bergmann P, Fauser M, Sattlegger D, Steger C (2019). "MVTec AD --- A Comprehensive Real-World Dataset for Unsupervised Anomaly Detection". In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). IEEE. pp. 9584–9592. ISBN 978-1-72813-293-8.
9. <sup>^</sup>Zou Y, Jeong J, Pemula L, Zhang D, Dabeer O (2022). "SPot-the-Difference Self-Supervised Pre-training for Anomaly Detection and Segmentation". In: European Conference on Computer Vision. p. 392–408. Available from: [arXiv:2207.14315](https://arxiv.org/abs/2207.14315).
10. <sup>^</sup>Chicco D, Jurman G (2020). "The Advantages of the Matthews Correlation Coefficient (MCC) over F1 Score and Accuracy in Binary Classification Evaluation". BMC Genomics. 21: 1–13.

11. <sup>a</sup> <sup>b</sup>Sokolova M, Japkowicz N, Szpakowicz S (2006). "Beyond Accuracy, F-Score and ROC: A Family of Discriminant Measures for Performance Evaluation". In: *AI 2006: Advances in Artificial Intelligence, Lecture Notes in Computer Science*. 4304: 1015–1021.
12. <sup>△</sup>Yi J, Yoon S (2020). "Patch SVDD: Patch-level SVDD for Anomaly Detection and Segmentation". In: *Asian Conference on Computer Vision (ACCV)*. arXiv. Available from: [arXiv:2006.16067](https://arxiv.org/abs/2006.16067).
13. <sup>△</sup>Defard T, Setkov A, Loesch A, Audigier R (2021). "Padim: A Patch Distribution Modeling Framework for Anomaly Detection and Localization". In: *International Conference on Pattern Recognition*. Springer. pp. 475–489.
14. <sup>△</sup>Gu Z, Zhu B, Zhu G, Chen Y, Tang M, Wang J (2024). "AnomalyGPT: Detecting Industrial Anomalies Using Large Vision-Language Models". In: *AAAI Conference on Artificial Intelligence*, vol. 38, pp. 1932–1940. arXiv:2308.15366.
15. <sup>△</sup>Roth K, Pemula L, Zepeda J, Schölkopf B, Brox T, Gehler P (2022). "Towards total recall in industrial anomaly detection." In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 14318–14328.
16. <sup>△</sup>Jeong J, Zou Y, Kim T, Zhang D, Ravichandran A, Dabeer O (2023). "WinCLIP: Zero-/Few-Shot Anomaly Classification and Segmentation". In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 19606–19616. Available from: [arXiv:2303.14814](https://arxiv.org/abs/2303.14814).
17. <sup>a</sup> <sup>b</sup>Radford A, Kim JW, Hallacy C, Ramesh A, Goh G, Agarwal S, Sastry G, Askell A, Mishkin P, Clark J, Krueger G, Sutskever I (2021). "Learning Transferable Visual Models From Natural Language Supervision". In: *International Conference on Machine Learning*. pp. 8748–8763.
18. <sup>△</sup>Liu H, Li C, Wu Q, Lee YJ (2023). "Visual Instruction Tuning". In: *Advances in Neural Information Processing Systems*, vol. 36. Available from: [arXiv:2304.08485](https://arxiv.org/abs/2304.08485).
19. <sup>a</sup> <sup>b</sup>Chen S, Han Z, He B, Buckley M, Torr P, Tresp V, Gu J (2024). "Understanding and Improving In-Context Learning on Vision-language Models". arXiv. Available from: [arXiv:2311.18021](https://arxiv.org/abs/2311.18021).
20. <sup>△</sup>Li B, Zhang Y, Chen L, Wang J, Yang J, Liu Z (2023). "Otter: A Multi-Modal Model with In-Context Instruction Tuning". arXiv. arXiv:2305.03726. Available from: <https://arxiv.org/abs/2305.03726>.
21. <sup>△</sup>Awadalla A, Gao I, Gardner J, Hessel J, Hanafy Y, Zhu W, Marathe K, Bitton Y, Gadre S, Sagawa S, Jitsev J, Kornblith S, Koh PW, Ilharco G, Wortsman M, Schmidt L (2023). "OpenFlamingo: An Open-Source Framework for Training Large Autoregressive Vision-Language Models". arXiv. arXiv:2308.01390. Available from: <https://arxiv.org/abs/2308.01390>.



22. <sup>△</sup>Li B, Zhang Y, Chen L, Wang J, Pu F, Yang J, Li C, Liu Z (2023). "MIMIC-IT: Multi-Modal In-Context Instruction Tuning". arXiv. arXiv:2306.05425. Available from: <https://arxiv.org/abs/2306.05425>.
23. <sup>△</sup>Tai Y, Fan W, Zhang Z, Zhu F, Zhao R, Liu Z (2023). "Link-Context Learning for Multimodal LLMs". arXiv. arXiv:2308.07891. Available from: <https://arxiv.org/abs/2308.07891>.
24. <sup>△</sup>Deng J, Dong W, Socher R, Li L-J, Li K, Fei-Fei L (2009). "Imagenet: A Large-Scale Hierarchical Image Database". In: IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 248–255. ISBN 1-4244-3992-2.
25. <sup>△</sup>Ba JL, Kiros JR, Hinton GE (2016). "Layer Normalization". arXiv. arXiv:1607.06450. Available from: <https://arxiv.org/abs/1607.06450>.
26. <sup>△</sup>Meta (2023). "Llama 2: Open Foundation and Fine-Tuned Chat Models". arXiv. arXiv:2307.09288. Available from: <https://arxiv.org/abs/2307.09288>.
27. <sup>△</sup>OpenAI (2023). "GPT-4 Technical Report". arXiv. arXiv:2303.08774. Available from: <https://arxiv.org/abs/2303.08774>.
28. <sup>△</sup>Baldassini FB, Shukor M, Cord M, Soulier L, Piwowarski B (2024). "What Makes Multimodal In-Context Learning Work?" arXiv. arXiv:2404.15736. Available from: <https://arxiv.org/abs/2404.15736>.
29. <sup>△</sup>Yang Z, Gan Z, Wang J, Hu X, Lu Y, Liu Z, Wang L (2022). "An Empirical Study of GPT-3 for Few-Shot Knowledge-Based VQA". In: AAAI Conference on Artificial Intelligence, vol. 36, series 3, pp. 3081–3089. arXiv:2109.05014.
30. <sup>△</sup>Steck H, Ekanadham C, Kallus N (2024). "Is Cosine-Similarity of Embeddings Really About Similarity?" In: Companion Proceedings of the ACM on Web Conference 2024. p. 887–890. Available from: arXiv [cs]. [2403.05440](https://arxiv.org/abs/2403.05440).
31. <sup>△</sup>He K, Zhang X, Ren S, Sun J (2015). "Deep Residual Learning for Image Recognition". In: IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 770–778. arXiv:1512.03385.
32. <sup>△</sup>{XTuner Contributors} (2023). "XTuner: A Toolkit for Efficiently Fine-tuning LLM."
33. <sup>△</sup>Loshchilov I, Hutter F (2019). "Decoupled Weight Decay Regularization". arXiv. arXiv:1711.05101. Available from: <https://arxiv.org/abs/1711.05101>.
34. <sup>△</sup>Grandini M, Bagli E, Visani G (2020). "Metrics for Multi-Class Classification: An Overview". arXiv. arXiv:2008.05756. Available from: <https://arxiv.org/abs/2008.05756>.
35. <sup>△</sup>G\u00f6gens M, Zhiyanov A, Tikhonov A, Prokhorenkova L (2022). "Good Classification Measures and How to Find Them". Neural Information Processing Systems. 34: 17136–17147. arXiv:2201.09044.

36. <sup>△</sup>Agarwal R, Singh A, Zhang LM, Bohnet B, Rosias L, Chan S, Zhang B, Anand A, Abbas Z, Nova A, Co-Reyes JD, Chu E, Behbahani F, Faust A, Larochelle H (2024). "Many-Shot In-Context Learning". arXiv. arXiv:2404.11018. Available from: <https://arxiv.org/abs/2404.11018>.
37. <sup>△</sup>Bertsch A, Ivgi M, Alon U, Berant J, Gormley MR, Neubig G (2024). "In-Context Learning with Long-Context Models: An In-Depth Exploration". arXiv. arXiv:2405.00200. Available from: <https://arxiv.org/abs/2405.00200>.
38. <sup>△</sup>Ren P, Xiao Y, Chang X, Huang P-Y, Li Z, Gupta BB, Chen X, Wang X (2021). "A Survey of Deep Active Learning". ACM computing surveys (CSUR). 54 (9): 1–40. arXiv:2009.00236.
39. <sup>△</sup>Ueno S, Yamada Y, Nakatsuka S, Kato K (2023). "Benchmarking of Query Strategies: Towards Future Deep Active Learning". arXiv. arXiv:2312.05751. Available from: <https://arxiv.org/abs/2312.05751>.
40. <sup>△</sup>Hu EJ, Shen Y, Wallis P, Allen-Zhu Z, Li Y, Wang S, Wang L, Chen W (2021). "LoRA: Low-Rank Adaptation of Large Language Models". arXiv. arXiv:2106.09685. Available from: <https://arxiv.org/abs/2106.09685>.
41. <sup>△</sup>Biderman D, Ortiz JG, Portes J, Paul M, Greengard P, Jennings C, King D, Havens S, Chiley V, Frankle J, Blakeney C, Cunningham JP (2024). "LoRA Learns Less and Forgets Less". arXiv. arXiv:2405.09673. Available from: <https://arxiv.org/abs/2405.09673>.
42. <sup>△</sup>Jiang X, Li J, Deng H, Liu Y, Gao B-B, Zhou Y, Li J, Wang C, Zheng F (2024). "Mmad: The first-ever comprehensive benchmark for multimodal large language models in industrial anomaly detection". arXiv preprint arXiv:2410.09453.

## Declarations

**Funding:** No specific funding was received for this work.

**Potential competing interests:** No potential competing interests to declare.