

Generalized Out-of-Distribution Detection and Beyond in Vision Language Model Era: A Survey

Atsuyuki Miyai, *Student Member, IEEE*, Jingkang Yang, Jingyang Zhang, Yifei Ming, Yueqian Lin, Qing Yu, *Member, IEEE*, Go Irie, *Member, IEEE*, Shafiq Joty, Yixuan Li, *Member, IEEE*, Hai Li, *Fellow, IEEE*, Ziwei Liu, *Member, IEEE*, Toshihiko Yamasaki, *Member, IEEE*, Kiyoharu Aizawa, *Fellow, IEEE*

Abstract—Detecting out-of-distribution (OOD) samples is crucial for ensuring the safety of machine learning systems and has shaped the field of OOD detection. Meanwhile, several other problems are closely related to OOD detection, including anomaly detection (AD), novelty detection (ND), open set recognition (OSR), and outlier detection (OD). To unify these problems, a generalized OOD detection framework was proposed, taxonomically categorizing these five problems. However, Vision Language Models (VLMs) such as CLIP have significantly changed the paradigm and blurred the boundaries between these fields, again confusing researchers. In this survey, we first present a generalized OOD detection v2, encapsulating the evolution of AD, ND, OSR, OOD detection, and OD in the VLM era. Our framework reveals that, with some field inactivity and integration, the demanding challenges have become OOD detection and AD. In addition, we also highlight the significant shift in the definition, problem settings, and benchmarks; we thus feature a comprehensive review of the methodology for OOD detection, including the discussion over other related tasks to clarify their relationship to OOD detection. Finally, we explore the advancements in the emerging Large Vision Language Model (LVLM) era, such as GPT-4V. We conclude this survey with open challenges and future directions.

Index Terms—Anomaly Detection, Novelty Detection, Open Set Recognition, Out-of-Distribution Detection, Outlier Detection, Vision Language Model, CLIP, Large Vision Language Model, Large Multi-modal Model

1 INTRODUCTION

A reliable visual recognition system should not only accurately predict known contexts, but also identify and reject unknown examples [1], [2], [3], [4]. In critical applications like autonomous driving, the system must alert and cede control to the driver upon encountering unfamiliar scenes or objects not seen during training. However, most existing machine learning models are trained based on the closed-world assumption [5], [6], where the test data is assumed to be drawn *i.i.d.* from the same distribution as the training data, known as in-distribution (ID). Therefore, the development of classifiers capable of detecting out-of-distribution (OOD) samples is a crucial challenge for real-world applications. This challenge is precisely the focus of

research in the field of OOD detection.

While OOD detection primarily focuses on semantic distribution shift, several other tasks share similar goals and motivations, including outlier detection (OD) [7], [8], [9], [10], anomaly detection (AD) [11], [12], [13], [14], novelty detection (ND) [15], [16], [17], [18], and open set recognition (OSR) [19], [20], [21]. Subtle differences in the specific definitions among these sub-topics have caused confusion in the field, leading to similar approaches being proposed across them. To address this issue, the generalized OOD detection framework was introduced [22]. The taxonomy of the generalized OOD detection framework is shown in Fig. 1. The generalized OOD detection framework classifies these tasks as special cases or sub-tasks under a unified taxonomy. This framework provides clear definitions and fosters a deeper understanding of each field.

In recent years, the emergence of Vision Language Models (VLMs), represented by CLIP [23], has rapidly accelerated research in the field of Computer Vision. This has changed the paradigm of the recognition field, allowing for zero-shot [23] or few-shot learning [24], [25] in various domains. VLMs have significantly influenced the aforementioned five problems (OD, AD, ND, OSR, and OOD detection), and the application of VLMs has become a highly notable research field [26], [27], [28], [29]. However, alongside this remarkable progress, the paradigm shift with the advent of the VLMs has blurred the boundaries between the five problems. Due to the difficulty of a clear understanding of the distinctions and interrelations between these tasks, each community within the fields is facing significant challenges in identifying the optimal direction to pursue in this VLM era.

In this survey, we introduce a novel unified framework

- A. Miyai, Q. Yu, T. Yamasaki, and K. Aizawa are with the Department of Information and Communication Engineering, The University of Tokyo, Japan, 1138656. Q. Yu is also with LY Corporation, Japan, 1028282. E-mail: {miyai,yamasaki}@com.t.u-tokyo.ac.jp, {yu,aizawa}@hal.t.u-tokyo.ac.jp
- J. Yang and Z. Liu are with S-Lab, School of Computer Science and Engineering, Nanyang Technological University, Singapore, 639798. E-mail: {jingkang001,ziwei.liu}@ntu.edu.sg
- J. Zhang, Y. Lin, and Hai Li are with the Department of Electrical and Computer Engineering, Duke University, Durham, NC, United States, 27708. E-mail: {jz288,yueqian.lin,hai.li}@duke.edu
- Y. Ming and S. Joty are with Salesforce AI Research, Palo Alto, CA, United States, 94301. S. Joty is also with NTU, Singapore, 639798 (on leave). E-mail: {yifei.ming, sjoty}@salesforce.com
- G. Irie is with the Department of Information and Computer Technology, Tokyo University of Science, Japan, 1258585. E-mail: goirie@ieee.org
- Y. Li is with the Department of Computer Sciences, University of Wisconsin-Madison, Madison, WI, United States, 53706. E-mail: sharonli@cs.wisc.edu

Resources: <https://github.com/AtsuMiyai/Awesome-OOD-VLM>.

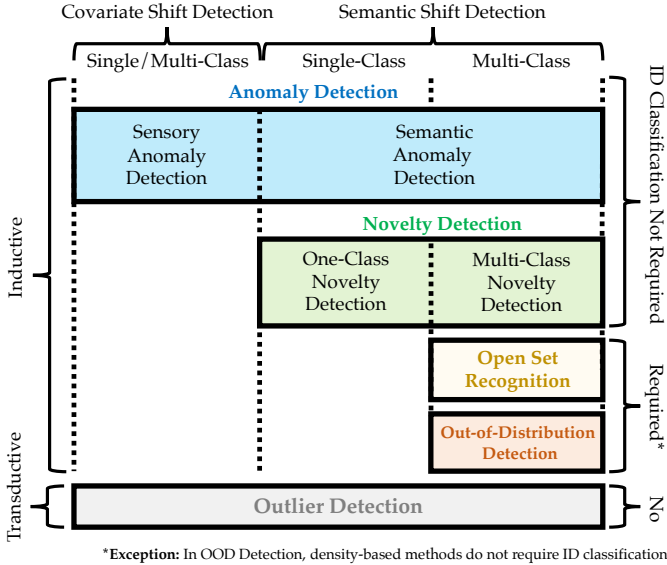


Fig. 1: Taxonomy of generalized OOD detection framework [22], illustrated by classification tasks. Figure adapted from [22].

termed *generalized OOD detection v2*, which extends the previous generalized OOD detection framework and summarizes the evolution of these five problems in the VLM era. To create it, we systematically review the use of VLMs across these five problem areas, tracing their development from the start to the present, and summarize the evolutionary trajectory of each problem. Importantly, our framework reveals that a paradigm shift has caused some fields to become inactive or integrate with others, and the demanding challenges in the VLM era become AD and OOD detection, which is a remarkable finding for each community. In addition to the inter-field evolution, we elaborate on the important shifts in the definition of OOD detection as well as the problem settings and benchmarks, with the contrast of those for related tasks. Then, we conduct a thorough review of the methodology for OOD detection and related tasks in the VLM era, intending to clarify their similarities and differences and inspire future research in OOD detection.

Finally, we introduce the evolution of these problems in the emerging Large Vision Language Model (LVLM) era, such as GPT-4V [30] or LLaVA [31] (also referred to as Large Multimodal Models or Multi-modal Large Language Models [32]). We summarize the definition of each evolving problem, the findings so far, and future challenges.

To sum up, in this survey paper, we make three contributions to the research community:

- 1) **Unified Framework in the VLM era:** We systematically review the evolution of five closely related topics of OD, AD, ND, OSR, and OOD detection in the Vision Language Model era and provide an updated unified framework termed *generalized OOD detection v2*. Our framework reveals that the paradigm shift has led to some field inactivity or integration, and the demanding challenges are AD and OOD detection. We hope that these observations highlight the demanding challenges in the VLM era and foster collaborative efforts among each

TABLE 1: Number of VLM-based papers in the Top Venues by June 2024.

Task	Top Venue
(a) Sensory AD	CVPR2023×1 [27] ICLR2024×1 [29] CVPR2024×4 [33], [34], [35], [36]
(b) Semantic AD/ND	TMLR2022×1 [37] (one-class) CVPR2024×1 [35] (one/multi-class)
(c) OSR	None
(d) OOD Detection	NeurIPS2021×1 [38] AAAI2022×1 [39] NeurIPS2022×1 [26] ICCV2023×1 [40] IJCV2023×1 [41] NeurIPS2023×3 [28], [42], [43] ICLR2024×2 [44], [45] CVPR2024×2 [46], [47] ICML2024×1 [48]
(e) OD	ICML2024 [49]

community.

- 2) **A Comprehensive Survey for OOD Detection in the VLM era:** While comprehensive surveys on OD, AD, ND, OSR, and OOD detection methodologies have been published in recent years [11], [12], [13], [14], [20], [22], [50], [51], this survey is the first to comprehensively overview OOD detection methods specifically in the VLM era. By connecting with other related tasks, we aim to provide readers with a holistic understanding of the developments and interconnections within each problem, particularly in the context of OOD detection.
- 3) **An Introduction to the Evolution in the LVLM Era:** We further introduce the evolution of each problem in the Large Vision Language Model era. Despite the infant stage of these fields, this survey offers an in-depth introduction to each problem, aiming to facilitate future advancements in this area.
- 4) **Future Research Directions:** We draw readers' attention to the future work necessary for advancing the field in the VLM and LVLM era. We conclude this survey with discussions on open challenges and opportunities for future research.

The paper content is organized as follows. In Sec. 2, we introduce the new version of generalized OOD detection by summarizing the evolution of the five related fields in the VLM era. We then overview the two key problems (OOD detection and AD) that have evolved and remain active in Sec. 3, with a detailed breakdown of existing methodologies being presented in Sec. 4 (CLIP-based OOD detection) and Sec. 5 (CLIP-based AD). In Sec. 6, we introduce early advancements of OOD detection and AD in the LVLM era. Sec. 7 and Sec. 8 feature discussions on potential challenges and future directions. Finally, we conclude with Sec. 9.

2 GENERALIZED OOD DETECTION V2

In this section, we introduce a novel unified framework termed *generalized OOD detection v2*, which summarizes the

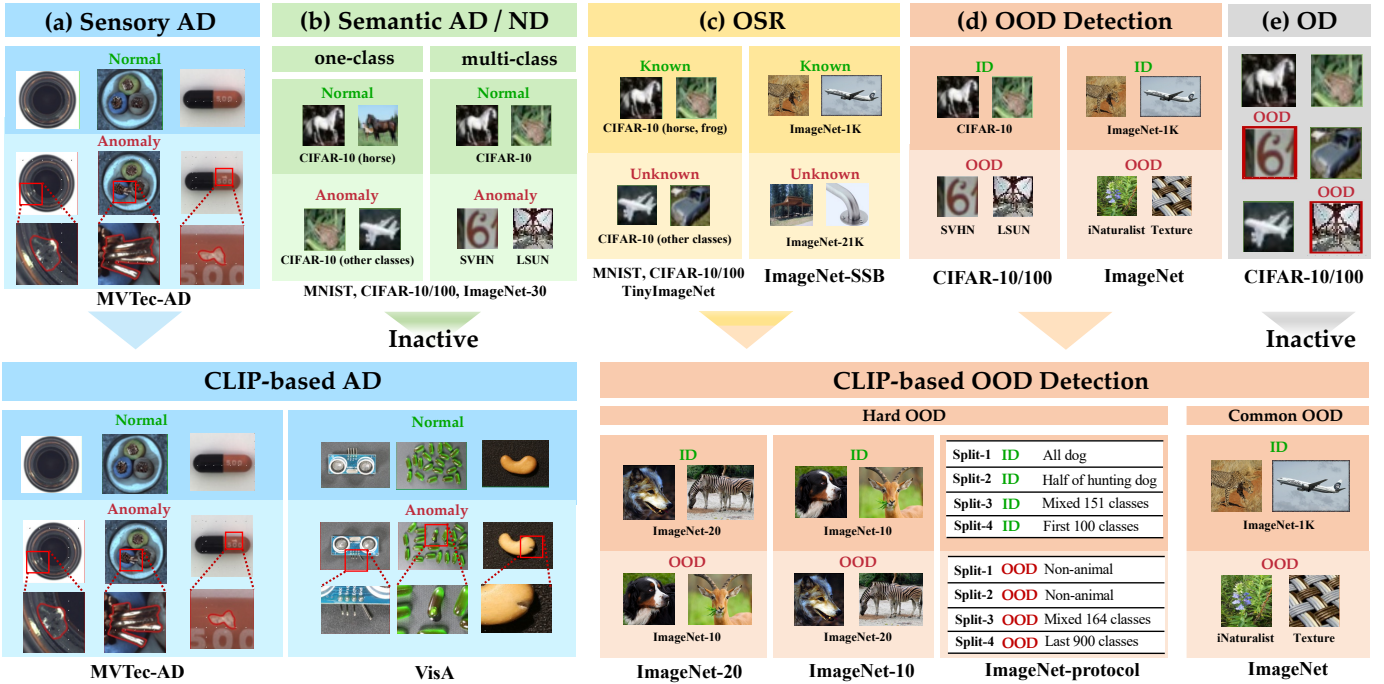


Fig. 2: Generalized OOD detection framework v2, reflecting the evolution of each problem in the VLM era. (a) Sensory AD has become a highly active and noteworthy field in the VLM era. In terms of benchmarks, in addition to the commonly used MVTec-AD [52], VisA [53], the largest industrial anomaly detection dataset, has also become a standard benchmark in the field. (b) Semantic AD/ND has become inactive in the VLM era. (c) OSR has been integrated into hard OOD detection. CLIP-based hard OOD detection incorporates the benchmark setup of OSR and creates new benchmarks such as ImageNet-10/ImageNet-20 [26] and ImageNet-protocol [47], [54]. (d) OOD detection is a highly active research area in the VLM era. (e) OD has become inactive in the VLM era.

evolution of the five related fields in the VLM era. We first revisit the previous generalized OOD detection framework in Sec. 2.1. Next, we introduce the evolution of each problem.

2.1 Background: Generalized OOD Detection V1

We first briefly revisit a previous *generalized OOD detection*, which encapsulates five related sub-topics: anomaly detection (AD), novelty detection (ND), open set recognition (OSR), out-of-distribution (OOD) detection, and outlier detection (OD). These sub-topics can be similar in the sense that they all define a certain *in-distribution*, with the common goal of detecting *out-of-distribution* samples under the open-world assumption. Previously, subtle differences existed among the sub-topics in terms of the specific definition and properties of in-distribution (ID) and OOD data.

To provide a clear definition, a generalized OOD detection framework was proposed [22]. The taxonomy for generalized OOD detection is shown in Fig. 1. It is based on the following four bases: (1) Distribution shift to detect: The task focuses on detecting either covariate shift (*e.g.*, OOD samples from a different domain) or semantic shift (*e.g.*, OOD samples from a different semantic). (2) ID data type: The in-distribution (ID) data contains either a single class or multiple classes. (3) Whether the task requires ID classification: Some tasks require classification of the ID data, while others do not. (4) Transductive vs. inductive learning: Transductive tasks require all observations (both ID and OOD), while inductive tasks follow the common train-test scheme. According to the above taxonomy, these five problems can be clearly

categorized as shown in Fig. 1: Anomaly detection is categorized into sensory anomaly detection, which deals with covariate shift, and semantic anomaly detection, which deals with semantic shift. Novelty detection falls under the same category as semantic anomaly detection. When addressing a multi-class scenario that necessitates ID classification, both open-set recognition and out-of-distribution detection are encompassed within this category. The main difference between OSR and OOD detection was the benchmark setup [22], [55] (Sec. 2.2 (c)). Outlier detection belongs to a different category from the other tasks, as this problem is transductive (*i.e.*, it has access to all observations).

For the detailed definition of each task, we refer the readers to the previous generalized OOD detection survey paper [22].

2.2 Evolution of Each Problem in VLM Era

We review how each problem has evolved in the VLM era. To make a fair judgment, we comprehensively investigated papers that use VLMs from top venues and summarized them in Table 1. Our survey revealed that CLIP [23] is predominantly used as the VLM for OOD detection and other sub-tasks, and other VLMs [56], [57] are rarely utilized. Therefore, we focus on CLIP as the target VLM in this survey and refer to OOD detection using CLIP as CLIP-based OOD detection. Similarly, we will prefix other tasks with “CLIP-based” (*e.g.*, CLIP-based AD). As OOD detection research is primarily focused on the image domain, we conduct a survey of other tasks within the image domain that are common

and have strong connections to OOD detection research. For instance, our survey does not cover video domain tasks [58], [59], [60] due to their limited connection to OOD detection.

(a) Sensory AD \rightarrow CLIP-based AD Sensory AD has continued to develop as a common problem setting for CLIP-based AD, inheriting the challenges of traditional sensory AD [27], [29], [33], [35], [61], [62], [62], [63], [64], [65]. As shown in Table 1, the first appearance in a top venue was at CVPR 2023, and since then, a total of six papers have been published in top venues. In addition, there are numerous other papers [61], [63], [64], [65], [66]. Moreover, in terms of benchmarks, in addition to the commonly used MVTec-AD [52], the largest industrial anomaly detection dataset VisA [53] has also become a standard benchmark in the field. Therefore, it is evident that sensory AD has become a highly active and noteworthy field in the VLM era.

(b) Semantic AD/ND \rightarrow Inactive Research on semantic AD/ND appears to become inactive in the VLM era. As shown in Table 1, there are only two papers, TMLR 2022 [37] and CVPR 2024 [35]. However, the CVPR 2024 work [35] aims to build a generalist anomaly detector that solves many AD tasks, including sensory AD and semantic AD, and is not primarily focused on semantic AD. The reasons for the inactivity include saturation of performance for one-class semantic AD/ND, and incompatibility of methods with CLIP for multi-class semantic AD/ND. As for one-class semantic AD/ND, TMLR [37] exists, but the performances with common CIFAR and ImageNet-30 datasets have already achieved around 99%. As for multi-class semantic AD/ND, a common approach is to treat ID classes as a single class, but treating ID classes as a single class is less compatible with CLIP’s class-wise discriminative capability.

(c) OSR \rightarrow CLIP-based OOD Detection We consider that OSR has been integrated into CLIP-based hard OOD detection. According to Table 1, there are no top venue publications on OSR research in the VLM era. Originally, the main difference between OSR and OOD detection was the benchmark setup [22], [55]. OSR typically divides the classes in the one dataset into some known (ID) classes and unknown (OOD) classes, as seen in MNIST-4/6 [67] CIFAR-4/6 [68], CIFAR-50/50 [69], and TinyImageNet-20/180 [70]. However, in recent years, some works on CLIP-based OOD detection incorporate the benchmark setup of OSR and create new benchmarks such as ImageNet-10/ImageNet-20 [26] and ImageNet-protocol [47], [54] for hard OOD detection. Therefore, the boundary between OOD detection and OSR has effectively disappeared, and all research in the VLM era has been integrated into OOD detection.

Nevertheless, while pure OSR research is declining, some studies have used the term “open-set” in the context of domain generalization [71]. These studies deviate from the original scope of OSR research and are rather closely aligned with the field of domain generalization [72]. Therefore, within our generalized OOD detection v2, we do not classify these studies as falling under OSR research. We will discuss them in the context of full-spectrum OOD detection, a research field combining generalization and detection, in Sec. 7.2.

(d) OOD Detection \rightarrow CLIP-based OOD Detection OOD detection is a highly active research area in the VLM era. As shown in Table 1, there are many papers in top venues,

indicating a high interest from the community. Additionally, as mentioned above, OSR has been integrated with OOD detection as a field of hard OOD detection [26], [47]. Therefore, it is expected that OOD detection will continue to grow and develop further.

(e) OD \rightarrow Inactive OD has become less active in the VLM era. Previously, OD was used for open-set semi-supervised learning [73], [74], [75], learning with open-set noisy labels [76], and novelty discovery [77], [78], [79], [80], [81]. The reason for the inactivity is that the use of CLIP led to a reduction in training costs and only a small amount of data needs to be collected, eliminating the need for large amounts of unlabeled data and reducing the need to consider noisy data. However, recently, Liang *et al.* [49] proposed Unsupervised Universal Fine-Tuning, a new problem setting for CLIP-based OD in ICML2024. Unsupervised Universal Fine-Tuning assumes a more realistic problem setting for unsupervised tuning of the downstream task with CLIP where some OOD samples are included in the unlabeled samples. With this new problem setting, there is still a possibility that OD will become active in the future. However, as OD is not currently an active area, we exclude OD from the main discussion of this survey. Unsupervised Universal Fine-Tuning is deeply related to OOD detection and will be discussed in detail in Sec. 4.3.

2.3 Discussion

Through Sec. 2.2, we found that previously mixed fields have been correctly organized in the VLM era, and that the focus should be on OOD detection and sensory AD. These fields are still developing, with an increasing number of methodologies and benchmarks, and are expected to become more active in the future. Note here that this does not mean that other fields have come to an end. For example, one reason why one-class semantic AD/ND has not been studied is the saturation of performance [37]. If more fine-grained and challenging datasets could be constructed, the field could be reactive. We put this in out-of-scope for this survey paper, but this is an important future challenge.

3 OVERVIEW OF EACH PROBLEM IN VLM ERA

In addition to the above inter-field evolution, we emphasize that the advent of VLMs has significantly changed the field of OOD detection itself. In this section, we present an overview of CLIP-based OOD detection, highlighting the key changes in the problem definition, the problem setting, and benchmarks. In addition, we also present an overview of CLIP-based AD in the hope that the understanding of each field will lead to a deeper understanding of CLIP-based OOD detection. For items that remain unchanged from traditional problems, such as background, applications, and evaluation, we refer the readers to the original generalized OOD detection paper [22].

3.1 CLIP-based Out-of-Distribution Detection

Definition The definition of CLIP-based OOD detection differs significantly from that of conventional OOD detection. Conventional OOD detection aims to detect test samples drawn from a distribution that is different from the training

distribution. As another definition, OOD detection is defined as a task to detect test samples that the model cannot or does not want to generalize [22]. However, for CLIP-based OOD detection, CLIP has a vast amount of knowledge, so the OOD sample is completely unrelated to the distribution of the CLIP’s pretraining data or the CLIP’s own generalization ability. Therefore, traditional definitions cannot adequately describe the definition of CLIP-based OOD detection.

Unlike the previous definition, CLIP-based OOD detection is defined as follows [26], [39]: CLIP-based OOD detection aims to detect samples that do not belong to any ID class text provided by the user. Given a pre-trained model, a classification task of interest is defined by a set of class labels \mathcal{Y}_{ID} , which we refer to as the ID classes. The semantic distribution is represented by the distribution $P(\mathcal{Y}_{ID})$. CLIP-based OOD detection aims to detect test samples that come from the distribution with the semantic shift from the ID classes, *i.e.*, $P(\mathcal{Y}_{ID}) \neq P(\mathcal{Y}_{OOD})$. Following the definition of the generalized OOD detection framework [22], ideal OOD detectors should keep the classification performance on test samples from ID class space \mathcal{Y}_{ID} , and reject OOD test samples with semantics outside the support of \mathcal{Y}_{ID} .

Problem Setting CLIP-based OOD detection focuses on solving the image classification task in a computationally efficient way. Unlike traditional OOD detection settings, which primarily involve training an ID classifier with whole ID data, CLIP-based OOD detection primarily focuses on a zero-shot [26] (*i.e.*, without utilizing ID images) or few-shot [28] (*i.e.*, utilizing only a few ID images) setting. Each detailed definition of both settings are described later in Sec. 4. The field is advancing towards greater computational efficiency, requiring minimal or no training data.

Benchmark Most recent works in CLIP-based OOD detection use high-resolution and large-scale datasets such as ImageNet [26], [28], [46], [47], [48]. The common ImageNet OOD benchmark uses ImageNet as ID and other datasets [82], [83], [84], [85] as OOD. However, in this common benchmark, the semantics between ID and OOD are far, which may allow easy distinction between the ID and OOD. Therefore, recent works use more challenging OOD benchmarks where they split ImageNet classes into ID and OOD categories for hard OOD detection [26], [47], [86]. The representative datasets are ImageNet-20 [26], ImageNet-10 [26], and the recently proposed ImageNet-protocol [54] created by dividing into multiple variations of ID/OOD pairs from ImageNet-1K. This creation strategy initially focused on OSR but has recently been repurposed for OOD detection. These changes in the datasets shift OOD detection closer to the real world and make it a more challenging and practical task.

3.2 CLIP-based Anomaly Detection

Definition Unlike OOD detection, the definition of anomaly detection (AD) has not changed between conventional AD and CLIP-based AD. AD is intended for use in specific circumstances (industrial inspection), where samples that deviate from predefined normality are considered an anomaly [11], [22]. Whether a model can generalize is irrelevant to the definition of “Anomaly”. Therefore, even with the emergence of CLIP, the definition has not changed.

Problem Setting CLIP-based AD focuses on solving anomaly classification and segmentation in a computationally efficient way. Anomaly classification, like conventional AD, is a binary classification task that distinguishes between normality and abnormality. Anomaly segmentation, also following conventional AD, involves segmenting the location of anomalies. Like CLIP-based OOD detection, CLIP-based AD also primarily focuses on a zero-shot [27] (*i.e.*, without utilizing images in the target dataset) or few-shot [27] (*i.e.*, using only a few normal images in the target dataset) setting. Each detailed definition of zero-shot and few-shot settings is described later in Sec. 5. As another shift, conventional AD created separate models for each category [87], [88], [89], [90], [91], [92], [93], while CLIP-based AD creates a single unified model across multiple categories [27], [29], [35], [61], [63], which leads to a more computationally efficient approach.

One key difference from CLIP-based OOD detection is that CLIP-based OOD detection does not involve localization tasks, while these are mainstream in CLIP-based AD. This will be discussed in detail in Sec. 5.4.

Benchmark Most works on CLIP-based AD tackle industrial inspection [52], [94], [95]. As for the benchmarks, besides the common MVTEC-AD dataset [52], the more challenging dataset VisA [53] has been newly employed [27]. The VisA benchmark includes objects with complex structures such as printed circuit boards and multiple instances with different locations within a single view, making it one of the most challenging datasets currently available in the open datasets. Since the pioneering work in CLIP-based AD (*i.e.*, WinCLIP [27]) used MVTEC-AD and VisA, many subsequent works have also used these datasets [33], [63], [65].

4 CLIP-BASED OOD DETECTION: METHODOLOGY

In this section, we introduce the methodologies for CLIP-based out-of-distribution (OOD) detection. Fig. 3 presents the timeline for representative methodologies for CLIP-based OOD detection. Table 2 presents representative methods. We introduce methods for zero-shot OOD detection in Sec. 4.1, few-shot OOD detection in Sec. 4.2, and other research directions in Sec. 4.3. For each methodology, we categorize them by the type of training and whether additional OOD prompts were employed.

4.1 Zero-shot Out-of-Distribution Detection

Zero-shot OOD detection was proposed in 2021 by Fort *et al.* [38]. Since then, a growing number of methods have been proposed year by year.

Definition of Zero-shot OOD Detection In zero-shot OOD detection, the term “Zero-shot” refers to the non-use of ID images during both training and inference phases. For instance, the method with additional training with auxiliary datasets (non-use of ID images) can be regarded as a zero-shot method [40]. The method with the pre-processing of the ID class texts can also be regarded as a zero-shot method [38], [39], [45], [48].

4.1.1 Training-free Methods

a. With OOD Prompts CLIP-based OOD detection started in this setting. The earliest work is ZeroOE [38]. ZeroOE

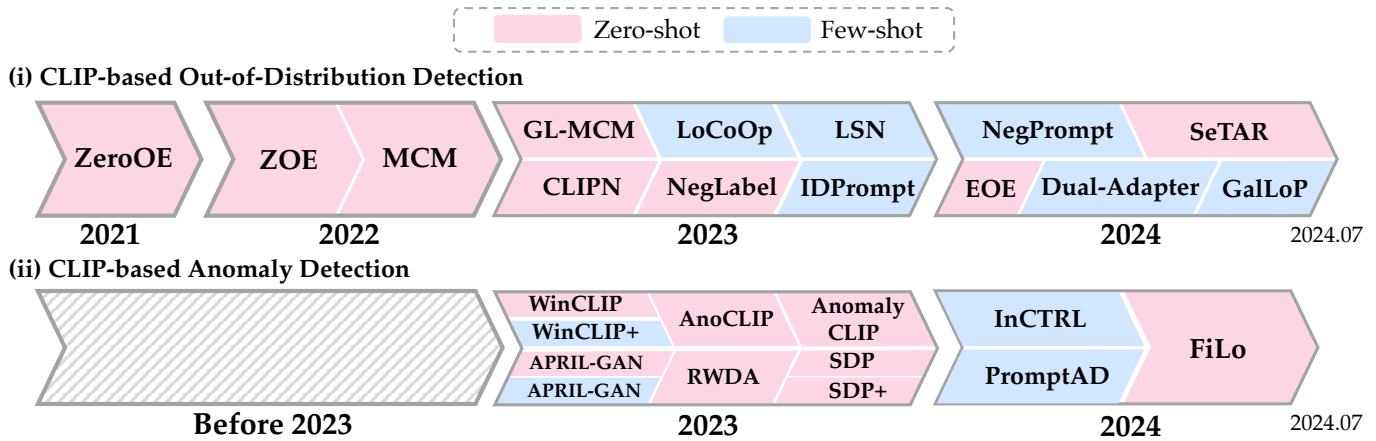


Fig. 3: Timeline for representative methodologies for CLIP-based out-of-distribution detection and CLIP-based anomaly detection. We observe that an increasing number of methods have recently been proposed for both tasks, indicating the growing activity in these fields.

TABLE 2: Representative paper list for CLIP-based out-of-distribution detection and CLIP-based anomaly detection.

Task	ID Image Availability	Training Type	OOD Prompts	Methods
§ 4 CLIP-based OOD Detection	§ 4.1 Zero-shot	§ 4.1.1 Training-free	✓	ZeroOE [38], ZOC [39], NegLabel [45], EOE [48]
		§ 4.1.2 Auxiliary Training	✗	MCM [26], GL-MCM [96], SeTAR [97]
	§ 4.2 Few-shot	§ 4.2.1 ID Training	✗	PEFT-MCM [41], LoCoOp [28], GalLoP [98]
		§ 4.2.2 Training-free	✓	LSN [44], NegPrompt [47], IDPrompt [46]
		§ 4.2.2 Auxiliary Training	✗	Dual-Adapter [99]
§ 5 CLIP-based Anomaly Detection	§ 5.1 Zero-shot	§ 5.1.1 Training-free	✓	WinCLIP [27], AnoCLIP [61], SDP [62]
		§ 5.1.2 Auxiliary Training	✓	APRIL-GAN (zero-shot) [63], RWDA [64], SDP+ [62], AnomalyCLIP [29], FiLo [65]
	§ 5.2 Few-shot	§ 5.2.1 Training-free	✓	WinCLIP+ [27]
		§ 5.2.2 ID Training	✓	PromptAD (one-class) [33]
		§ 5.2.3 Auxiliary Training + ref.	✓	APRIL-GAN (few-shot) [63], InCTRL [35]

feeds the potential OOD labels to the textual encoder of CLIP. However, the method of using known OOD labels is infeasible for real-world applications. To solve this problem, ZOC [39] proposed to train an OOD label generator based on the visual encoder of CLIP and use the generated pseudo-OOD labels for OOD detection. However, when dealing with large-scale datasets encompassing a multitude of ID classes, the label generator may not generate effective candidate OOD labels, resulting in poor performance. Building on these early works [38], [39], recent works focus on how to obtain high-quality OOD labels through either (i) OOD label retrieval [45], [100] or (2) OOD label generation [48]. (i) One of the representative retrieval-based methods is NegLabel [45]. NegLabel selects high-quality OOD labels from extensive corpus databases by calculating the distance between an extracted OOD label and ID label. (ii) One of the representative generation-based methods is EOE [48]. EOE utilizes Large Language Models (LLMs) to produce high-quality OOD labels. By modifying the prompts given to the LLM, EOE can be generalized to a variety of tasks, including far and near OOD detection.

b. Without OOD Prompts In zero-shot OOD detection,

many methods utilize OOD labels, but the difficulty and cost of creating these labels pose challenges. To address these issues, Ming *et al.* [26] proposed MCM, which uses only ID labels to detect OOD. MCM is a simple approach that devises softmax scaling to align visual features with textual concepts for OOD detection. Despite its simplicity, MCM has high effectiveness and scalability, and it serves as a crucial baseline in CLIP-based OOD detection. Building on the concept of MCM, Miyai *et al.* [96] proposed GL-MCM, which extends MCM by just adding a local MCM score to enhance the fine-grained detection capability in local regions. SeTAR [97] enhances MCM and GL-MCM by changing the model’s weight matrices using a simple greedy search algorithm. We consider these methods to be post-hoc methods for CLIP-based OOD detection in that they directly employ an ID classifier for OOD detection. Due to their simplicity and high scalability, these post-hoc methods can bring fundamental performance improvements for many subsequent methods [28], [44], [47]. Therefore, we expect that this field should be developed further in the future, reflecting the trajectory of the field before CLIP emerged [101], [102], [103], [104], [105], [106], [107], [108], [109], [110], [111].

4.1.2 Auxiliary Training-based Methods

CLIPN [40] is the only auxiliary training-based method for zero-shot OOD detection. CLIPN aims to empower the logic of saying “no” within CLIP, and it designs a novel learnable “no” prompt and an additional “no” text encoder to capture negation semantics within images. To create an additional text encoder, CLIPN needs to be pre-trained on the CC-3M dataset [112]. While the extensive pre-training of CLIPN may indeed lead to intensive computations and lower scalability, once pre-trained, it performs zero-shot OOD detection across a wide range of domains with comparable performance to few-shot OOD detection methods [28], [44]. In the future, within this field, there is potential for zero-shot open-vocabulary OOD detection to further advance the realm of zero-shot OOD detection, which will be discussed later in Sec. 7.2.

4.2 Few-shot Out-of-Distribution Detection

Few-shot OOD detection was concurrently proposed by Miyai *et al.* [28] and Ming *et al.* [41] in June 2023. Since then, it has become the most active research area in CLIP-based OOD detection.

Definition of Few-shot OOD Detection CLIP-based few-shot OOD detection aims to detect OOD images using only a few labeled ID images. In few-shot OOD detection, the term “Few-shot” refers to the use of a few ID images during training or inference phases. For instance, the method with additional training with a few ID images can be regarded as a few-shot method [33]. Even without training, if a method uses a few ID images as a reference, we regard it as a few-shot method [27]. Regarding the number of shots, it is common to experiment with 1-shot to 16-shot [28], [99], following the closed-set setting [24].

4.2.1 ID Training-based Methods

a. Without OOD Prompts Few-shot OOD detection began in this setting. Ming *et al.* [41] proposed PEFT-MCM for CLIP-based OOD detection, which demonstrates the effectiveness of combining parameter-efficient tuning methods (*e.g.*, prompt learning [24] or adapter [113]) and MCM [26]. Concurrently, Miyai *et al.* [28] proposed LoCoOp, a pioneer prompt learning approach for few-shot OOD detection. LoCoOp enhances CoOp’s [24] OOD detection capabilities by performing OOD regularization with local OOD features. LoCoOp is the simplest prompt learning method and serves as a crucial baseline in few-shot OOD detection. Unlike LoCoOp, which utilizes non-ID local regions for OOD regularization, GalLoP [98] proposes an approach that utilizes local ID regions to enable a more fine-grained distinction between ID and OOD samples. GalLoP learns a diverse set of prompts by utilizing both global and local visual representations, thereby enhancing the detection capabilities.

b. With OOD Prompts Similar to zero-shot OOD detection, recent works in few-shot OOD detection utilize additional OOD prompts [44], [46], [47]. As representative methods, LSN [44] and NegPrompt [47] were proposed concurrently. They state that the simple negative prompts added “not” (*e.g.*, “not a photo of a [cls]”) fail to capture the dissimilarity for identifying OOD samples. Therefore, by preparing negative

prompts and training with them, LSN and NegPrompt can learn suitable negative prompts, enabling more accurate detection of OOD samples. The difference between LSN and NegPrompt lies in their approach to the use of negative prompts. LSN prepares unique negative prompts for each class and learns suitable negative prompts for each class. In contrast, NegPrompt prepares multiple negative prompts common to all ID classes and trains them to learn generic templates representing the negative semantics of any given class labels. Also, NegPrompt tested the performances in the hard OOD detection setting with ImageNet-protocol [54], outperforming LoCoOp and CoOp. Alternatively, IDPrompt [46] takes a different approach by introducing ID-like prompts, which are designed for OOD features that are close to the ID features. It extracts ID-like OOD regions in ID training images and trains ID-like prompts with these extracted OOD data. In a unique direction, LAPT [114] proposes an automatic sample collection strategy that retrieves or generates training ID images only with ID class names, which achieves high performance without image collection and annotation costs. LAPT then performs distribution-aware prompt learning, which distinguishes between ID class and OOD class tokens. LAPT is positioned within the context of more efficient few-shot OOD detection in this survey paper since it requires generating or retrieving “ID images” for the data collection.

In the context of few-shot OOD detection, recently, Li *et al.* [47] proposed a new problem setting called open-vocabulary OOD (OV-OOO) detection. While common few-shot OOD detection involves training on images from all ID classes during training, OV-OOO detection involves training on images from just a small subset of ID classes and performing OOD detection using all ID classes at inference time. Formally, we define a subset of semantic labels $\mathcal{Y}_{ID,sub} \subset \mathcal{Y}_{ID}$, where \mathcal{Y}_{ID} represents all ID labels. Based on this subset of labels, we define a corresponding subset dataset $\mathcal{D}_{ID,sub}^{train} \subset \mathcal{D}_{ID}^{train}$. During training, only $\mathcal{D}_{ID,sub}^{train}$ is used. Then, at inference time, all ID classes \mathcal{Y}_{ID} are used, and the goal is to detect OOD from a combination of all ID test data \mathcal{D}_{ID}^{test} with \mathcal{Y}_{ID} and OOD test data \mathcal{D}_{OOD}^{test} with \mathcal{Y}_{OOD} . For this setting, existing few-shot OOD detection methods [28], [47] can be easily applied by simply combining the rest of the ID classes. In particular, NegPrompt [47] learns general negative prompts that are not specific to the training ID classes, so it achieves high performance in OV-OOO detection.

4.2.2 Training-free Methods

Training-free few-shot OOD detection is a novel research field, and only Dual-Adapter [99] falls under this category. Dual-Adapter adopts a prior-based method Tip-Adapter [113], which leverages both textual and visual features with a cache model and enhances performance without training. To adapt this to few-shot OOD detection, Dual-Adapter employs the concept of dual cache modeling and constructs Positive-Adapter and Negative-Adapter, and identifies OOD samples with the prediction difference with both adapters.

4.3 Other Important Research Directions

4.3.1 CLIP-based Full-spectrum OOD Detection

CLIP-based full-spectrum OOD (FS-OOD) detection is a crucial challenge [115]. FS-OOD detection was proposed by Yang *et al.* [116] in 2022 as an important setting that considers both OOD generalization [117], [118] and OOD detection simultaneously. Unlike standard OOD detection, which only focuses on semantic shifts between training and test distributions, FS-OOD detection further considers non-semantic covariate shift by including covariate-shifted ID images. As for the benchmarks, OpenOOD v1.5 [119] provides two large-scale benchmarks based on ImageNet-200 and ImageNet-1K, incorporating ImageNet-C [117] with image corruptions, ImageNet-R [118] with style changes, and ImageNet-V2 [120] with resampling bias as ID. As for CLIP-based methods, LSA [115] uses a bidirectional prompt customization mechanism, which adjusts discriminative ID and OOD boundary.

4.3.2 Other Tasks with CLIP-based OOD Detection

Unsupervised Universal Fine-Tuning CLIP-based OOD detection is useful for a new task called Unsupervised Universal Fine-Tuning (UUFT) [49]. UUFT is a problem of unsupervised learning for outlier detection (OD). Existing studies for unsupervised learning assumed that all unlabeled images belong to one of the ID classes [121], [122], [123], but they require prior knowledge of exact class names linked to ground truth labels, which restricts their usefulness in various real-world situations. For a more realistic setting, UUFT assumes that OOD images are included in the unlabeled images. To detect OOD images during training, they developed MCM [26] and proposed UEO, which leverages sample-level confidence to approximately minimize the conditional entropy of confident instances and maximize the marginal entropy of less confident instances.

Open-world Prompt Tuning CLIP-based OOD detection is useful for a new task called Open-world Prompt Tuning [124]. Open-world Prompt Tuning is a task that evaluates the classification accuracy on a mix of known and novel ID classes while training the model with known classes. To solve this problem, Zhou *et al.* [124] proposed DeCoOp which incorporates OOD detection into the inference pipelines and improves the base-to-new separability, preventing performance degradation on new classes.

5 CLIP-BASED AD: METHODOLOGY

In this section, we introduce methodologies for CLIP-based anomaly detection (AD) in the hope that the contrast with OOD detection clarifies the similarities and differences between each task and facilitates a deeper understanding of CLIP-based OOD detection.

5.1 Zero-shot Anomaly Detection

CLIP-based zero-shot AD was proposed in 2023 by Jeong *et al.* [27]. Although it started about two years later than OOD detection, many methods have been proposed up to the present.

Definition of Zero-shot AD The meaning of the term “Zero-shot” for zero-shot AD is similar to that for zero-shot

OOD detection. In zero-shot AD, the term “Zero-shot” refers to the non-use of the images in the target domain during both training and inference phases. For instance, the method with additional training with auxiliary datasets can be regarded as a zero-shot method [29], [62], [63], [64], [65]. The method with the pre-processing of the target class texts can also be regarded as a zero-shot method [27], [61], [62].

5.1.1 Training-free Methods

With Anomaly Prompts In zero-shot AD, a common approach is to utilize anomaly prompts to detect anomalies. This hypothesis is supported by several observations from existing work [27]. Firstly, the concepts of normality and anomalies are context-dependent states [125] of an object, with language playing a crucial role in defining these states. Secondly, language provides additional insights that help differentiate defects from acceptable variations in normality.

The simplest zero-shot AD methods are (i) to perform anomaly classification with CLIP using text prompts for normality and anomalies as classes (*i.e.*, “normal [class]” *vs.* “anomalous [class]”) and (ii) to calculate the similarity to the normal prompt (*i.e.*, “normal [class]”) as the score. These methods are called CLIP-AC [27]. Jeong *et al.* [27] reported that CLIP-AC with both normal and anomaly prompts outperforms that with only normal text prompts, which indicates the importance of the use of anomaly prompts. However, the performances for this naive method are not yet satisfactory due to the wide range of variations of anomalies. To solve this issue, Jeong *et al.* [27] proposed WinCLIP. WinCLIP performs a compositional ensemble on a large number of pre-defined normal and anomaly templates and efficient extraction and aggregation of window/patch/image-level features aligned with the text. WinCLIP outperforms CLIP-AC by a large margin. Because of its simplicity and pioneering work, WinCLIP has become an important baseline for CLIP-based AD. AnoCLIP [61] follows WinCLIP’s approach of using a large number of pre-defined normal and anomaly templates but modifies the templates to be domain-aware (*e.g.*, industrial photo) and contrastive state for normal and anomaly (*e.g.*, perfect and imperfect). However, it is noteworthy that the performance of the ensemble strategies of previous methods heavily depends on the text descriptions [27], [61]. Also, it is observed that more descriptions are not always better [62], which makes the previous approaches [27], [61] using a naive ensemble of large templates somewhat uncontrollable and random in their applications. Therefore, SDP [62] proposes RVS, a representative vector selection paradigm, which makes the mechanism of extracting representative vectors from large templates controllable, allowing for a more diverse selection of representative vectors.

5.1.2 Auxiliary Training-based Methods

Unlike CLIP-based OOD detection, all auxiliary training-based methods in zero-shot AD are open-vocabulary AD methods that are trained with auxiliary AD datasets and tested on unseen target datasets by simply changing the category prompt. Existing methods perform supervised training on the test set of one dataset and perform zero-shot testing on the other dataset [29], [63], [65] (*e.g.*, training with MVTEC-AD and evaluation with VisA.)

In recent years, the development of auxiliary training-based methods has received more attention than training-free zero-shot methods. There are two main reasons why training is necessary for AD: (i) The first is the domain gap between semantics and anomalies. CLIP is pre-trained to understand the semantics of images, so when applied in a zero-shot manner, it captures the semantics of the image. However, actual anomalies are not semantics, but rather represent the state of an object and appear only in local areas of the image. Therefore, without training, this domain gap between semantics and anomalies cannot be bridged. (ii) The second reason is that there are limitations to relying on a large set of manually crafted anomaly prompts. This incurs prompt creation costs and also makes it difficult to respond to unknown anomalies. Therefore, by replacing the anomaly prompts with learnable parameters, they aim to solve the high costs and limited adaptability to new anomalies.

To address the above issue (i), APRIL-GAN (also known as VAND) [63] was proposed. APRIL-GAN tackles the domain gap between semantics and anomaly by adding additional linear layers in vision encoders. These linear layers project image features at each scale into the text space, creating and aggregating anomaly maps at each stage. Similarly, SDP+ [99] also incorporates additional linear layers into SDP [99] to effectively project image features into the text feature space, addressing the misalignment between image and text. To solve both the issue (i) and (ii), AnomalyCLIP [29] was proposed. AnomalyCLIP is a prompt learning-based method similar to CoOp. By replacing anomaly prompts with learnable parameters, it eliminates the need to prepare a large number of manually pre-defined prompts such as those in WinCLIP [27]. Furthermore, unlike CoOp which learns object semantics, AnomalyCLIP learns object-agnostic text prompts that capture generic normality and abnormality in an image regardless of its semantics. To achieve this, AnomalyCLIP introduces object-agnostic text prompt templates for both normal and anomaly and performs global and local context optimization. A more recent approach, FiLo [65] leverages Large Language Models (LLMs) to generate fine-grained anomaly descriptions for each object category. This method replaces generic abnormal descriptions with LLM-generated specific anomaly content for each sample. By adding learnable prompts before the generated anomaly prompts, FiLo performs global and local context optimization, enhancing the ability to detect anomalies. As a unique direction from these methods, RWDA [64] proposes a data augmentation approach by utilizing CLIP’s text embeddings as training data. RWDA adds randomly generated words into normal and anomaly prompts to generate a diverse set of normal and anomaly training samples and trains a regular feed-forward neural network with diverse text embeddings.

5.2 Few-shot Anomaly Detection

CLIP-based few-shot AD was proposed in 2023 by Jeong *et al.* [27], concurrently with the development of zero-shot AD [27]. Traditional few-shot AD research focuses on modeling the normal distribution from a limited number of normal samples to detect anomalies [126], [127], [128], [129], [130], [131], [132], [133]. However, these methods often struggle

to generalize to new domains, as they typically require retraining with the target datasets. With the advent of CLIP, the field of few-shot AD is shifting towards using a few target images only for inference at test time without training.

Definition of Few-shot AD CLIP-based few-shot AD aims to detect anomaly images using only a few images in the target domain. The meaning of the term “Few-shot” is similar to that of few-shot OOD detection. In few-shot AD, the term “Few-shot” refers to the use of a few normal images in the target domain during training or inference phases. For instance, the method with additional training with a few normal images in the target domain can be regarded as a few-shot method [28], [41]. Even without training, if a method uses a few normal images in the target domain as a reference, we regard it as a few-shot method [27].

5.2.1 Training-free Methods

The earliest approach in CLIP-based few-shot AD is WinCLIP+ [27], an improved method of WinCLIP. WinCLIP, a base zero-shot AD method, cannot identify certain defects that can only be defined visually rather than textually. For example, the “Metal-nut” category in MVTecAD has an anomaly type labeled “flipped upside down,” which can only be identified relative to a normal image. To address this, WinCLIP+ incorporates a few normal reference images into a memory bank [134] and calculates the anomaly score with the cosine similarity between the query image and its most similar image in the memory bank.

5.2.2 ID Training-based Methods

There are few training-based methods for CLIP-based few-shot AD research, although CLIP-based few-shot OOD detection actively explores the training-based methods. This is likely because detecting anomalies in unknown classes is highly valuable in practical applications, and training on target data may oversimplify the task of few-shot anomaly detection. This oversimplification concern arises because the anomaly space for known categories is much more limited compared to that for OOD detection. Training with the target data could therefore make the task oversimplified, reducing its difficulty. The only existing work in this area is PromptAD [33], a prompt learning method for one-class AD (where the normal class consists of one class). In one-class AD, traditional prompt learning methods for multi-class classification (*e.g.*, CoOp [24]) do not work well. To address this, PromptAD creates a large number of anomaly prompts by adding a learnable anomaly suffix to the normal prompt. It then learns to bring the visual features closer to the normal prompt and further away from the anomaly prompts, enabling prompt learning for one-class AD.

5.2.3 Auxiliary Training- and Reference-based Methods

We explore the methods trained on auxiliary datasets and utilize the normal images in the target domain as references during inference. An early work in this category is APRIL-GAN (few-shot) [63], which uses a linear layer trained with auxiliary datasets. Similar to WinCLIP+ [27], APRIL-GAN (few-shot) utilizes a few ID reference images with a memory bank-based approach [134]. More recently, Zhu *et al.* [35] proposed an in-context-learning-based method called

InCTRL. InCTRL trains a model to discriminate anomalies from normal samples by learning to identify residuals or discrepancies between query images and a set of few-shot normal images (in-context sample prompts) from auxiliary data. During inference, InCTRL identifies anomalies by measuring the discrepancy between the features of the query image and a few in-context normal samples from the target dataset.

5.3 Other Research Direction

5.3.1 Anomaly Detection with Localization Models

Some works [135], [136] tackle AD using SAM [57] or DINO [56], foundation models for localization. The representative work is SAA and SAA+ [135]. SAA is a simple baseline approach and utilizes Grounding-DINO [56] for the anomaly region generator and SAM [57] for the anomaly region refiner. SAA+ is an improved method of SAA and it incorporates domain expert knowledge and target image context into SAA. Given that AD necessitates localization, it is expected that the number of works employing localization foundation models such as SAM will continue to increase.

5.3.2 Medical Anomaly Detection

While most works on CLIP-based AD focus on industrial AD, recent studies have begun to challenge medical anomaly detection (medical AD) [35], [36], [62], [137], [138]. CLIP-based medical AD is a more challenging area than industrial AD due to the larger gap between different data modalities. A representative work on medical AD is MVFA [36]. MVFA is a method specifically tailored for medical AD. It incorporates multiple residual adapters into the CLIP’s visual encoder to reduce the domain gap, enabling a stepwise enhancement of visual features across different levels. The future progression of medical AD and industrial AD offers an intriguing perspective, exploring whether these fields will develop independently or influence each other. However, when considering practical applications, it should be noted that medical AD faces the challenge that anomalies are not always describable in the language. Therefore, the development of medical AD methods that do not use CLIP is also important.

5.4 Discussion

We discuss the similarities and differences between CLIP-based OOD detection and CLIP-based AD to deepen our understanding of CLIP-based OOD detection.

5.4.1 Difference between Each Methodology

Differing Scopes of OOD OOD detection and AD differ significantly in the scope of OOD (anomaly) they cover, which leads to differences in methodologies, particularly in the use of OOD prompts. As explained in Sec. 3, sensory AD is intended for specific use cases like industrial inspection, where samples deviating from predefined normality (e.g., defective products) are considered anomalies [11], [22]. In other words, in sensory AD, the anomaly space is limited to damaged objects with shared semantics, and anomalies like images of dogs are not expected. This limited anomaly space allows even simple prompts to achieve decent performance. Therefore, as shown in Table 2, all AD methods

utilize anomaly prompts. Conversely, in OOD detection, as explained in Sec. 3.1, anything semantically different from the ID class is considered OOD. Thus, utilizing naive manual OOD prompts is forbidden (even if it improves benchmark performance). This vastness of the OOD space is the key factor differentiating the methodologies between the two fields.

Essential Features Learned for Detecting OOD There is a significant difference in the features that need to be captured between CLIP-based OOD detection and CLIP-based AD. In CLIP-based OOD detection, it is ideal to learn a more compact ID decision boundary that produces low uncertainty for the ID data, with high uncertainty for OOD data [28]. On the other hand, CLIP-based AD aims to learn the anomalies, instead of ID decision boundary for detecting anomalies [29]. The learned features are entirely different, which is helpful in developing each training-based method.

Difficulty of Localization Task CLIP-based OOD detection and CLIP-based AD differ significantly in the difficulty of OOD (anomaly) localization tasks. In CLIP-based AD, anomaly segmentation is a mainstream task, often performed alongside classification in many papers. However, in CLIP-based OOD detection, there has been no research on object-level OOD detection/segmentation. Object-level OOD detection aims to detect OOD objects [139], [140], [141]. The inactivity is related to the size of the OOD space, and the too-vast space of OOD makes it difficult to identify OOD objects with prompts effectively. To pave the way for future development, the foundation models for localization such as SAM [57], which can segment individual objects, have the potential to address object-level OOD detection/segmentation. Object-level OOD detection/segmentation using SAM is a promising future research direction.

5.4.2 Similarity between Each Methodology

Each Problem Setting The existing problem settings for CLIP-based AD and CLIP-based OOD detection are similar. Both primarily focus on zero-shot and few-shot settings and can be categorized into training-free, auxiliary training-based, and ID training-based methods. By examining each problem setting more closely, we can observe, for instance, that while open-vocabulary AD is predominant in CLIP-based AD, it has not been deeply explored in CLIP-based OOD detection. This provides valuable insights into future directions that are crucial for OOD detection.

History of Approaches The history of the progress of methods for CLIP-based AD and CLIP-based OOD detection is similar. For instance, both problems initially started with naive methods with manual OOD prompts (ZeroOE [38] for OOD detection, WinCLIP [27] for AD). To address the issues with these initial approaches, subsequent methods emerged that replaced OOD prompts with learnable parameters (LSN [44] and NegPrompt [47] for OOD detection, and AnomalyCLIP [29] for AD). Therefore, by carefully examining each other’s fields, there is potential for mutual enhancement and interaction in the future.

6 EVOLUTION IN LVL M ERA

In this section, we introduce the early advances in OOD detection and AD in the Large Vision Language Models

(LVLM) era. While previous sections focused on VLMs such as CLIP, this section shifts our focus to the more emerging topic of “Large” VLMs. Recent advancements in computer vision have led to the emergence of LVLMs such as GPT-4V [30] and LLaVA [31]. Although these fields are still in their early stages with limited papers, this survey provides a deep introduction to each problem in the hope that our detailed review can help foster further advancements in this area.

6.1 Change of Each Problem

i. Sensory AD → Sensory AD Sensory AD has continued to develop in the LVLM era [142], [144], [145]. The use of LVLMs has made AD applicable in many domains and modalities [142].

ii. OOD Detection → Unsolvable Problem Detection In the LVLM era, OOD detection has evolved into a new task termed Unsolvable Problem Detection (UPD) [143]. UPD evaluates the LVLMs’ ability to recognize and abstain from answering unexpected or unsolvable input questions, effectively expanding the scope of OOD detection into the context of Visual Question Answering (VQA) tasks. This shift to the VQA task has significantly broadened the concept of OOD detection to a wider range of AI tasks involving LVLMs.

6.2 Unsolvable Problem Detection

6.2.1 Summary of Problem

Background Following the recent revolutionary development of LLMs [146], [147], [148], [149], [150], [151], LVLMs [149], [152], [153], [154], [155], [156], [157], [158], [159] have demonstrated remarkable capabilities in diverse applications [160], [161], [162], [163]. However, a significant concern has arisen regarding the reliability of these models, specifically their ability to generate accurate and trustworthy information. These models frequently produce incorrect or misleading information, a phenomenon referred to as “hallucination” [32]. Among the various hallucination issues [32], the challenge of identifying out-of-place questions is crucial for deploying LVLMs in safety-critical applications. This challenge extends the concept of OOD detection to the VQA tasks for LVLMs and represents a specific aspect of LVLMs’ trustworthiness.

Definition Unsolvable Problem Detection (UPD) is a task to measure the trustworthiness of LVLMs, which is designed to evaluate models’ capacity to withhold answers when faced with unsolvable problems. The UPD task can be categorized into three distinct problem types: Absent Answer Detection (AAD), Incompatible Answer Set Detection (IASD), and Incompatible Visual Question Detection (IVQD). The details of each setting are as follows:

- 1) **Absent Answer Detection (AAD):** AAD evaluates the model’s capacity to determine when the correct answer is absent from the provided options.
- 2) **Incompatible Answer Set Detection (IASD):** IASD assesses the model’s ability to discern answer choices that are completely irrelevant to the given question and image.

- 3) **Incompatible Visual Question Detection (IVQD):** IVQD evaluates the model’s capacity to discern whether a question and image are unrelated or mismatched.

Benchmark Miyai *et al.* [143] created MM-UPD Bench for the UPD challenge. MM-UPD encompasses MM-AAD, MM-IASD, and MM-IVQD benchmarks for each UPD problem. Each benchmarks are created on the top of MMBench (dev) [162], which is a systematically designed objective benchmark for evaluating various abilities of LVLMs. Following the definition of each ability in MMBench (*e.g.*, “Coarse Perception: Image Scene” and “Logic Reasoning: Future Prediction”), MM-UPD evaluates the trustworthiness of LVLMs from various abilities.

Although MM-UPD is the main benchmark, the adaptation cost of creating UPD problems is not high, making it highly applicable to other benchmarks. For instance, the recently proposed MuirBench [164], a benchmark for multi-image understanding, has incorporated the concept of UPD by adding unsolvable problems.

Application UPD has a wide range of applications, from everyday use of LVLMs to robot manipulation. Especially when incorporating LVLMs into safety-critical domains such as robot manipulation [165] and autonomous driving [166], there is a risk of significant problems if the LVLM fails to identify erroneous user questions and makes incorrect predictions. UPD serves as a task to ensure safety in such safety-critical scenarios.

Evaluation UPD introduces new evaluation metrics that incorporate the concept of the evaluation protocols for OOD detection, taking into account the prediction distribution for both standard (ID) and UPD (OOD) samples. The rationale is that ideal LVLMs should not only give correct answers for the standard problems but need to withhold answering in the UPD scenario where the problem is unsolvable. To better reflect the ideal behavior of LVLMs, UPD measures several metrics: (i) **Standard Accuracy:** The accuracy on standard problems where the image, question, and answer sets are all aligned, and the ground-truth answer is always present within the provided options. (ii) **UPD Accuracy:** The accuracy of AAD/IASD/IVQD problems. (iii) **Dual Accuracy:** Accuracy of standard and UPD pairs. We count success only if the model is correct for both standard and UPD problems.

6.2.2 Findings

In the following, we briefly summarized the findings of the UPD challenge [143].

1. Most LVLMs Hardly Hesitate to Answer. Most LVLMs, especially open-source LVLMs, have significantly low UPD accuracies, which indicates the difficulty of the UPD challenge. For example, LLaVA-1.5 [31] and CogVLM [155], which are state-of-the-art LVLMs, completely fail to withhold answering. GPT-4V achieves higher performances than other LVLMs due to its safety training process [167]. However, there is still a performance gap from the upper bound scores.

2. Performance Tendency Differs a lot by Each Ability in the Benchmark. The performance of LVLMs differs in each ability in the MM-UPD Bench. For instance, GPT-4V has its

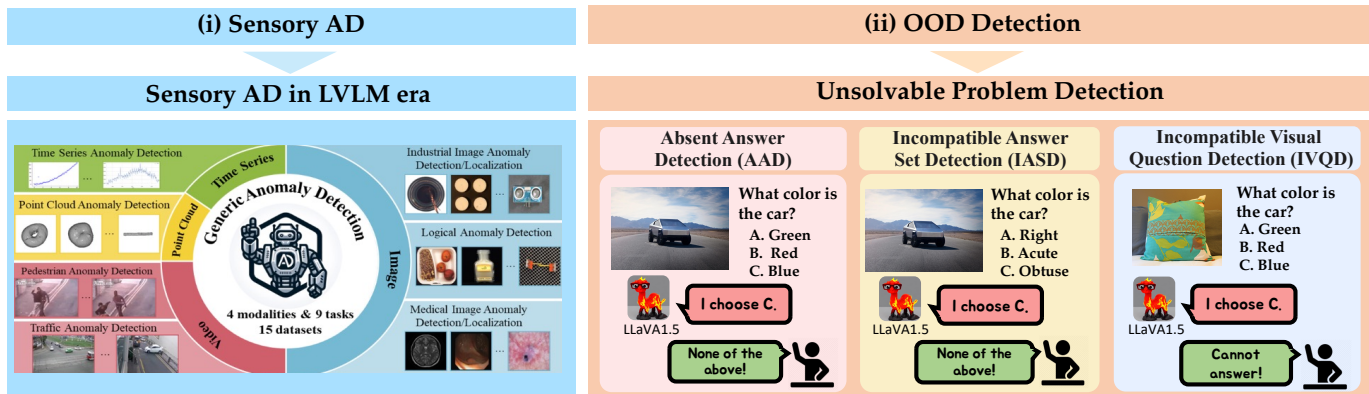


Fig. 4: Overview of the evolution of each problem in the Large Vision Language Model (LVLM) era. (i) Sensory AD has become a highly active and noteworthy field in the LVLM era [142]. (ii) Notably, OOD detection is evolving into a new task called Unsolvable Problem Detection [143] in the LVLM era. Figure adapted from [142] and partially from [143].

limitation in attribute comparison and LLaVA-NeXT-34B has its limitation in object localization.

3. Effective Prompt Strategies Vary Across Different LVLMs

Effective prompt strategies vary across different LVLMs. In the original paper, they experimented with an option-based prompt approach that adds an option of “None of the above” and an instruction-based approach that adds an instruction “If all options are incorrect, answer None of the above”. As a result, the effectiveness of each approach differs significantly depending on the type of LVLMs. This highlights the difficulty of finding an effective prompt strategy for all LVLMs.

6.3 Anomaly Detection in LVLM Era

6.3.1 Summary of Problem

Background Anomaly detection is a crucial task in a variety of domains and data types. However, existing anomaly detection models are often designed for specific domains or modalities [142]. Also, current AD methods only provide an anomaly score for the test sample and require a manual threshold to distinguish between normal and anomalous instances for each sample [144]. To facilitate real-world applications, developing a system capable of expressing anomalies in natural language across various modalities and domains is crucial for ensuring accessibility to a wider range of users.

Definition The definition of AD remains consistent with conventional and CLIP-based AD, aiming to identify samples that deviate from predefined normality. The key distinction lies in the output. While previous methods produced an anomaly score requiring a manual threshold, AD with LVLMs aims to recognize and describe anomalies using text, eliminating the need for manual thresholds and enhancing human interpretability.

Benchmark Since the field of AD with LVLMs is in its infant stage, there are still no unified benchmarks. AnomalyGPT [144] focuses on industrial image anomaly detection/localization and uses the standard benchmarks MVTEC-AD [52] and VisA [53]. More recently, Cao *et al.* [142] extend the domain and modality and demonstrate the applications in industrial image anomaly detection/localization (*e.g.*,

MVTec-AD [52]), point cloud anomaly detection (MVTec 3D [168]), medical image anomaly detection/localization (*e.g.*, Chest X-ray [169], Head CT [170]), logical anomaly detection (*e.g.*, MVTEC LOCO [171]), pedestrian anomaly detection (*e.g.*, UCF-Crime Dataset [172]), traffic anomaly detection (*e.g.*, Kaggle Accident Detection [173]), and time series anomaly detection (*e.g.*, Outlier Detection Dataset [174]).

Evaluation Evaluation in anomaly detection with LVLMs is an open challenge. AnomalyGPT [144] asks LVLMs the question “Is there an anomaly in this image?” and determines anomaly or normal based on the simple rule-based approach of whether the response contains a “yes” or “no”. However, this rule-based approach is not robust, as a response is considered correct even if the explanation following “yes” is completely incorrect. On the other hand, Cao *et al.* [142] conducted only qualitative evaluations and left quantitative evaluations as an open challenge. Therefore, the evaluation of anomaly detection by LVLM is a future challenge.

6.3.2 Findings

Cao *et al.* [142] described the observations of GPT-4V in the paper, so we briefly summarized them here.

1. GPT-4V Excels in Zero/One-shot Settings across Various Modalities and Fields. GPT-4V shows proficiency in identifying anomalies in multi-modality (*e.g.*, images, point clouds, X-rays) and multi-field (*e.g.*, industrial, medical, pedestrian, traffic, and time series anomaly detection). In addition, GPT-4V demonstrates strong performance in both zero-shot and one-shot settings.

2. GPT-4V can Understand Both Global and Fine-grained Anomalies. GPT-4V can recognize both global and local abnormal patterns or behaviors, which indicates the ability to understand global and fine-grained semantics.

3. GPT-4V can be Enhanced with Increasing Prompts. By giving more context and information, the model significantly improves its ability to detect anomalies accurately.

7 POTENTIAL CHALLENGES

This section discusses potential challenges for CLIP-based OOD detection that may be highlighted by the widespread adoption of our framework. Since similar or ambiguous

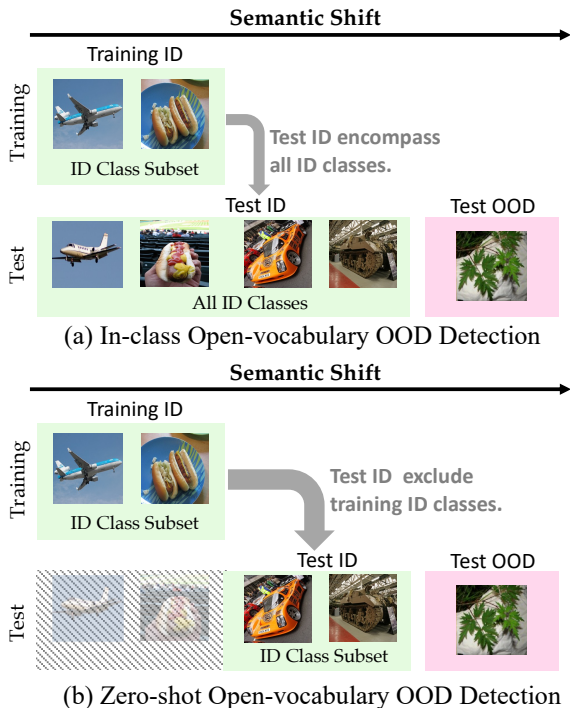


Fig. 5: Illustration of the settings for in-class open-vocabulary OOD detection and zero-shot open-vocabulary OOD detection. In-class open-vocabulary OOD detection should be compared with few-shot OOD detection methods, while zero-shot open-vocabulary OOD detection should be discussed in the context of zero-shot OOD detection.

problem settings currently exist, we present these challenges to avoid future confusion for readers.

7.1 Open-vocabulary OOD Detection

The existing setting for open-vocabulary OOD (OV-OOD) detection [47] involves training with images from a small subset of ID classes and evaluating with all ID class names, as explained in Sec. 4.2.1. However, this problem definition has several issues: (i) The OOD detection performance on ID classes not included in the training subset is unclear, (ii) This contradicts the settings of the existing open-vocabulary settings, including CLIP-based AD (Sec. 5.1.2), which assume no explicit class overlap between training and test ID data [29], [63], [65], [175], [176], leading to confusion in the fields. Therefore, for clear evaluation and common understanding in the field, we redefine the OV-OOD detection as in-class open-vocabulary OOD detection (in-class OV-OOD detection) and zero-shot open-vocabulary OOD detection (zero-shot OV-OOD detection). The illustration for each setting is shown in Fig. 5. Detailed explanations for each setting are as follows:

In-class OV-OOD Detection In-class OV-OOD detection is the same as the original setting proposed by [47], involving training with images from a small subset of ID classes and evaluating on all ID class names including training classes. In-class OV-OOD detection should be evaluated in the few-shot OOD detection setting to assess how well these methods retain their performance when trained with limited class data. In this setting, it is preferable to employ the existing

benchmark setup that randomly selects subset ID classes from all ID classes [47].

Zero-shot OV-OOD Detection Zero-shot OV-OOD detection shares the training setting with in-class evaluation, using a small subset of ID classes. However, it differs in the testing phase, where it exclusively utilizes classes not included in the training subset. This problem setting is consistent with the common open-vocabulary setting [175], [176] and the open-vocabulary setting in CLIP-AD [29], [63], [65]. Zero-shot OV-OOD detection can be positioned within the area of zero-shot OOD detection, which would facilitate the research of classes for training, we should ensure semantic dissimilarity between the classes in the evaluation and training sets to guarantee the validity of the zero-shot problem setting. For instance, when splitting ImageNet classes into training and evaluation subsets for zero-shot OV-OOD detection, it would be preferable to consider the hierarchical classes of ImageNet to avoid similar semantics between training and evaluation subsets.

7.2 CLIP-based Full-spectrum OOD Detection

CLIP-based full-spectrum OOD detection is a well-established task with existing research [115]. However, with the future prevalence of hard OOD detection, which incorporates the OSR setup, the distinction between this task and open-set domain generalization (OSDG) [177], [178], [179], [180] may become ambiguous. OSDG shares a similar motivation of realizing models capable of both generalization and detection. The illustration of both problem settings is shown in Fig. 6 (a)(b). OSDG assumes that only data with covariate shifts will be input during testing, and aims to classify input samples into one of ID classes if their semantics align with ID classes and detect OOD that exhibit semantic shifts. OSDG is a research area with a much smaller number of works compared to the main topics of this survey (OD, AD, ND, OSR, and OOD detection), but, in recent years, some works have tackled OSDG using CLIP [71], [181], [182].

To eliminate potential ambiguity, this paper defines a hard full-spectrum OOD (FS-OOD) detection, following our generalized OOD detection v2. The illustration of the problem setting for hard FS-OOD detection is shown in Fig. 6 (c). Hard FS-OOD detection expands the scope of existing FS-OOD detection by utilizing the conventional OSR setup. It introduces OOD samples from different categories within the training domain and OOD samples that share the same covariate shift as the target domain but exhibit semantic shifts. This task can be regarded as an extensive version of both existing FS-OOD detection and OSDG tasks. Hard FS-OOD detection represents a developing challenge with promising potential for future research.

8 FUTURE DIRECTIONS

In this section, we discuss the future directions of OOD detection and UPD. For OOD detection, we explore not only OOD detection for VLMs but also single-modal OOD detection, with a specific focus on emerging challenges as VLMs evolve. For a discussion of the long-standing challenges in OOD detection, we can refer the readers to the previous generalized OOD detection paper [22].

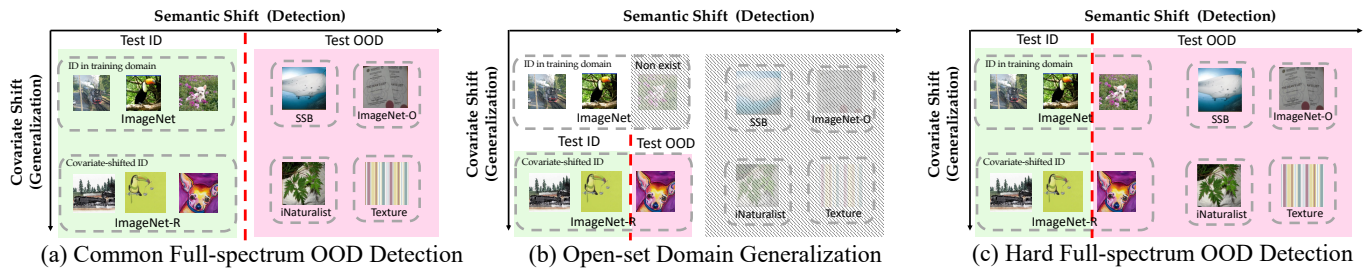


Fig. 6: Illustration of the settings for full-spectrum OOD detection, open-set domain generalization, and hard full-spectrum OOD detection. Hard full-spectrum OOD detection is a potential challenge arising from our generalized OOD detection v2.

8.1 OOD Detectin for Vision Language Models

a. Hard OOD Detection Hard OOD detection will become increasingly important in the future due to its high practicality and the challenging nature of the problem. Hard OOD detection utilizes the OSR benchmark setup, where some classes within a single dataset are designated as ID and others as OOD. In this field, not only small datasets such as ImageNet-10 and ImageNet-20 [26] but also datasets with a larger number of classes such as ImageNet-protocol [54] have been proposed. Many existing studies, such as LoCoOp [28] and LSN [44], primarily use the common ImageNet OOD benchmark, so hard OOD detection has not yet been well studied. This field will develop further in the future.

b. Post-hoc Methods To propose post-hoc methods is important for the fundamental performance improvement of CLIP-based OOD detection. The methods of directly employing an ID classifier such as MCM [26] are called post-hoc methods. Prior to CLIP, various approaches were proposed [101], [102], [103], [104], [105], [108], [109], [111], [183], [184]. However, CLIP-based post-hoc methods often underperform methods with additional OOD prompts, so they are not extensively researched in zero-shot OOD detection. However, we should focus on the scalability of post-hoc methods. The post-hoc methods [26], [96] can be easily applied to many subsequent methods [28], [44], [47], [185], bringing fundamental performance improvements. Furthermore, very recently, post-hoc methods specifically tailored for prompt learning methods have also emerged [86]. Therefore, proposing post-hoc methods and demonstrating the improvements not only in zero-shot but also in subsequent few-shot settings [28], [185] is crucial, even if they underperform methods using OOD prompts in the zero-shot setting. This field should continue evolving, mirroring its growth before the advent of CLIP.

c. Bridging the Gap with Closed-set Classifiers. OOD detection ensures the safety of ID classifiers, so it is crucial to bridge the gap between the advancements in existing closed-set classifiers and OOD detection. Currently, the representative method for few-shot OOD detection is LoCoOp [28], a text prompt learning method based on CoOp. However, in closed-set settings, few-shot learning methods based on text prompt learning have been proposed other than CoOp [186], [187], [188]. Furthermore, text-based prompt learning methods only train the text prompts, so they cannot handle differences in image domains. Therefore, adopting methods that can handle image domain differences [189], [190], [191] for OOD detection is essential for bridging the

gap with closed-set ID classifiers.

d. Training-free Few-shot OOD Detection The research direction of training-free few-shot OOD detection is still in its infant stage, with only one existing study [99]. Given the prevalence of training-based methods in few-shot OOD detection, proposing methods that do not require training is crucial. Considering the advancements of training-free methods in CLIP-based anomaly detection, we anticipate a similar trajectory for CLIP-based OOD detection. Future directions include refining adapter-based methods or leveraging external knowledge such as retrieval augmentation [192], [193]. Addressing training-free few-shot OOD detection is a pivotal step towards realizing more computationally efficient OOD detection in the future.

e. Full-spectrum OOD Detection CLIP-based full-spectrum OOD (FS-OOD) detection is a promising research area [115], [116]. In practical applications, there is a strong motivation to create models that can not only detect semantically shifted OOD inputs but also generalize to covariate-shifted data [116], [194]. Within CLIP-based methods, OOD detection and generalization are often discussed in separate contexts [28], [195], resulting in a trade-off between detection and generalization performance [98]. Furthermore, there exists potential ambiguity in the distinction between FS-OOD detection and open-set domain generalization. To eliminate the potential ambiguity, we have formulated a potential problem setting called hard FS-OOD detection. We hope this survey inspires further advancements and developments for FS-OOD detection.

f. Open-vocabulary OOD Detection Open-vocabulary OOD (OV-OOD) detection has a high practical potential, but it is still in its infant stages [47]. In particular, as explained in Sec. 7.1, zero-shot OV-OOD detection is a potential research field. We hope that this survey paper will inspire future efforts in OV-OOD detection.

g. Object-level OOD Detection Object-level OOD detection remains an unexplored area. As discussed in Sec. 5.4, this is due to the vastness of the OOD space, which makes it difficult to identify OOD objects using texts. To pave the way for future advancements in object-level OOD detection/segmentation, foundation models for localization, such as SAM [57], offer a promising solution. By integrating these models with methods such as MCM [26], we can potentially achieve object-level OOD detection and segmentation, opening up a new frontier in OOD detection research.

8.2 Single-modal OOD Detection

a. Leveraging Large Pre-trained Models Leveraging large pre-trained models is important for single-modal OOD detection. Numerous methods for OOD detection conduct experiments using backbones trained from scratch and do not utilize pre-trained models [22], [101], [104], [108], [109], [119], [196], [197]. In a recent study, Miyai *et al.* [198] systematically investigated the impact of pre-training on OOD detection from both the perspectives of the types of OOD data and pre-training algorithms [199], [200]. Dong *et al.* [201] explored parameter-efficient learning for single-modal OOD detection and proposed DSGF, which leverages both fine-tuned features and original pre-trained features. While leveraging large pre-trained models [202] with lightweight tuning is an active area of research in single-modal closed-set classification [191], [203], there have been limited studies for single-modal OOD detection, which presents a promising avenue for future research.

b. Real-world Benchmarks and Evaluations Considering the future development of CLIP-based OOD detection, there should be increasing focus on expanding the scope of benchmarks to encompass real-world scenarios where CLIP is less applicable. For instance, recently, Baek *et al.* [204] introduced ImageNet-ES, consisting of variations in environmental and camera sensor factors by directly capturing 202k images with an actual camera in a controlled testbed, which bridges the gap between the common benchmark and the real-world scenario. Besides, utilizing datasets such as WILDS [205], [206], which consider real-world data shifts, or datasets for medical OOD detection [207], can provide valuable insights, especially in safety-critical applications like autonomous driving and medical image analysis.

8.3 Unsolvable Problem Detection

a. Exploring Effective Solutions It is important to propose effective solutions for UPD. One of the potential approaches is to adapt the methodologies of OOD detection to UPD. For example, the perplexity of the LVLM’s response could be used as a score to identify unsolvable queries. Furthermore, proposing model-agnostic post-hoc methods is also important to enhance the reliability of many LVLMs. Therefore, incorporating the concepts of OOD detection techniques for UPD is an important direction for future work.

b. Extension to Diverse Benchmarks MM-UPD Bench consists of general QA datasets. However, UPD can be incorporated into more diverse benchmarks including domain-specific knowledge for advanced reasoning [163], [208] and multi-image understanding [164]. For example, Muir-Bench [164] incorporated the concept of UPD as a metric for robust evaluation. Integrating the concept of UPD into benchmarks is essential for evaluating the robustness and trustworthiness of LVLMs in their target tasks.

c. Theoretical Understanding of UPD Theoretically understanding the reasons behind the difficulty of UPD could provide the community with valuable insights. Theorizing the behavior of LVLMs poses a shared challenge in the field, highlighting the importance of collaborative efforts within the community.

9 CONCLUSION

In this survey, we comprehensively review the evolution of the five problems including AD, ND, OSR, OOD detection, and OD in the VLM era, and propose a framework of *generalized OOD detection v2*. Our framework identifies OOD detection and AD as the primary challenges in the VLM era, which highlights the demanding challenges and fosters collaborative efforts among each community. By articulating the shifts in the definitions, problem settings, and benchmarks, we encourage subsequent works to accurately understand their evolving target problems in the VLM era. By sorting out the methodologies, we hope that readers can easily grasp the mainstream methods, identify important baselines and novel problem settings, and propose future solutions. By shedding light on recent studies in the LVLM era, we hope that researchers within each community can identify promising research directions in this emerging era. By providing future directions, we hope that our survey will clarify the tasks to be tackled by future works in the VLM and LVLM era, thereby facilitating future advances in the right direction.

ACKNOWLEDGMENT

This work was supported by JST BOOST, Japan Grant Number JPMJBS2418 and JST JPMJCR22U4. We thank Toyooka Mashiro (AYM Lab at UTokyo) for his valuable assistance in designing the figures, and Kazuki Egashira, Yuki Imajuku, Takubon Son, and Zaiying Zhao (AYM Lab at UTokyo) for valuable feedback on the paper.

REFERENCES

- [1] D. Amodei, C. Olah, J. Steinhardt, P. Christiano, J. Schulman, and D. Mané, “Concrete problems in AI safety,” *arXiv preprint arXiv:1606.06565*, 2016. 1
- [2] S. Mohseni, H. Wang, Z. Yu, C. Xiao, Z. Wang, and J. Yadawa, “Practical machine learning safety: A survey and primer,” *arXiv preprint arXiv:2106.04823*, 2021. 1
- [3] D. Hendrycks, N. Carlini, J. Schulman, and J. Steinhardt, “Unsolved problems in ML safety,” *arXiv preprint arXiv:2109.13916*, 2021. 1
- [4] D. Hendrycks and M. Mazeika, “X-risk analysis for AI research,” *arXiv preprint arXiv:2206.05862*, 2022. 1
- [5] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *NIPS*, 2012. 1
- [6] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification,” in *ICCV*, 2015. 1
- [7] C. C. Aggarwal and P. S. Yu, “Outlier detection for high dimensional data,” in *ACM SIGMOD*, 2001. 1
- [8] V. Hodge and J. Austin, “A survey of outlier detection methodologies,” *Artificial intelligence review*, 2004. 1
- [9] I. Ben-Gal, “Outlier detection,” in *Data Mining and Knowledge Discovery Handbook*, 2005. 1
- [10] H. Wang, M. J. Bah, and M. Hammad, “Progress in outlier detection techniques: A survey,” *IEEE Access*, vol. 7, pp. 107964–108000, 2019. 1
- [11] L. Ruff, J. R. Kauffmann, R. A. Vandermeulen, G. Montavon, W. Samek, M. Kloft, T. G. Dietterich, and K.-R. Müller, “A unifying review of deep and shallow anomaly detection,” *Proceedings of the IEEE*, vol. 109, no. 5, pp. 756–795, 2021. 1, 2, 5, 10
- [12] G. Pang, C. Shen, L. Cao, and A. V. D. Hengel, “Deep learning for anomaly detection: A review,” *ACM Comput. Surv.*, vol. 54, no. 2, 2021. 1, 2
- [13] S. Bulusu, B. Kailkhura, B. Li, P. K. Varshney, and D. Song, “Anomalous example detection in deep learning: A survey,” *IEEE Access*, vol. 8, pp. 132330–132347, 2020. 1, 2

- [14] R. Chalapathy and S. Chawla, "Deep learning for anomaly detection: A survey," *arXiv preprint arXiv:1901.03407*, 2019. [1](#), [2](#)
- [15] M. A. Pimentel, D. A. Clifton, L. Clifton, and L. Tarassenko, "A review of novelty detection," *Signal processing*, vol. 99, pp. 215–249, 2014. [1](#)
- [16] D. Miljković, "Review of novelty detection methods," in *MIPRO*, 2010. [1](#)
- [17] M. Markou and S. Singh, "Novelty detection: a review—part 1: statistical approaches," *Signal processing*, vol. 83, no. 12, pp. 2481–2497, 2003. [1](#)
- [18] M. Markou and S. Singh, "Novelty detection: a review—part 2:: neural network based approaches," *Signal processing*, vol. 83, no. 12, pp. 2499–2521, 2003. [1](#)
- [19] T. E. Boulton, S. Cruz, A. R. Dhamija, M. Gunther, J. Henrydoss, and W. J. Scheirer, "Learning and the unknown: Surveying steps toward open world recognition," in *AAAI*, 2019. [1](#)
- [20] C. Geng, S.-j. Huang, and S. Chen, "Recent advances in open set recognition: A survey," *IEEE TPAMI*, vol. 43, no. 10, pp. 3614–3631, 2020. [1](#), [2](#)
- [21] A. Mahdavi and M. Carvalho, "A survey on open set recognition," *arXiv preprint arXiv:2109.00893*, 2021. [1](#)
- [22] J. Yang, K. Zhou, Y. Li, and Z. Liu, "Generalized out-of-distribution detection: A survey," *IJCV*, pp. 1–28, 2024. [1](#), [2](#), [3](#), [4](#), [5](#), [10](#), [13](#), [15](#)
- [23] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, *et al.*, "Learning transferable visual models from natural language supervision," in *ICML*, 2021. [1](#), [3](#)
- [24] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, "Learning to prompt for vision-language models," *IJCV*, vol. 130, no. 9, pp. 2337–2348, 2022. [1](#), [7](#), [9](#)
- [25] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, "Conditional prompt learning for vision-language models," in *CVPR*, 2022. [1](#)
- [26] Y. Ming, Z. Cai, J. Gu, Y. Sun, W. Li, and Y. Li, "Delving into out-of-distribution detection with vision-language representations," in *NeurIPS*, 2022. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [14](#)
- [27] J. Jeong, Y. Zou, T. Kim, D. Zhang, A. Ravichandran, and O. Dabeer, "Winclip: Zero-/few-shot anomaly classification and segmentation," in *CVPR*, 2023. [1](#), [2](#), [4](#), [5](#), [6](#), [7](#), [8](#), [9](#), [10](#)
- [28] A. Miyai, Q. Yu, G. Irie, and K. Aizawa, "Locoop: Few-shot out-of-distribution detection via prompt learning," in *NeurIPS*, 2023. [1](#), [2](#), [5](#), [6](#), [7](#), [9](#), [10](#), [14](#)
- [29] Q. Zhou, G. Pang, Y. Tian, S. He, and J. Chen, "Anomalyclip: Object-agnostic prompt learning for zero-shot anomaly detection," in *ICLR*, 2024. [1](#), [2](#), [4](#), [5](#), [6](#), [8](#), [9](#), [10](#), [13](#)
- [30] C. Lu, C. Qian, G. Zheng, H. Fan, H. Gao, *et al.*, "From gpt-4 to gemini and beyond: Assessing the landscape of mlms on generalizability, trustworthiness and causality through four modalities," *arXiv preprint arXiv:2401.15071*, 2024. [2](#), [11](#)
- [31] H. Liu, C. Li, Y. Li, and Y. J. Lee, "Improved baselines with visual instruction tuning," in *CVPR*, 2024. [2](#), [11](#)
- [32] Z. Bai, P. Wang, T. Xiao, T. He, Z. Han, Z. Zhang, and M. Z. Shou, "Hallucination of multimodal large language models: A survey," *arXiv preprint arXiv:2404.18930*, 2024. [2](#), [11](#)
- [33] X. Li, Z. Zhang, X. Tan, C. Chen, Y. Qu, Y. Xie, and L. Ma, "Promptad: Learning prompts with only normal samples for few-shot anomaly detection," in *CVPR*, 2024. [2](#), [4](#), [5](#), [6](#), [7](#), [9](#)
- [34] C.-H. Ho, K.-C. Peng, and N. Vasconcelos, "Long-tailed anomaly detection with learnable class names," in *CVPR*, 2024. [2](#)
- [35] J. Zhu and G. Pang, "Toward generalist anomaly detection via in-context residual learning with few-shot sample prompts," in *CVPR*, 2024. [2](#), [4](#), [5](#), [6](#), [9](#), [10](#)
- [36] C. Huang, A. Jiang, J. Feng, Y. Zhang, X. Wang, and Y. Wang, "Adapting visual-language models for generalizable anomaly detection in medical images," in *CVPR*, 2024. [2](#), [10](#)
- [37] P. Liznerski, L. Ruff, R. A. Vandermeulen, B. J. Franks, K.-R. Müller, and M. Kloft, "Exposing outlier exposure: What can be learned from few, one, and zero outlier images," *Transactions on Machine Learning Research*, 2022. [2](#), [4](#)
- [38] S. Fort, J. Ren, and B. Lakshminarayanan, "Exploring the limits of out-of-distribution detection," in *NeurIPS*, 2021. [2](#), [5](#), [6](#), [10](#)
- [39] S. Esmaeilpour, B. Liu, E. Robertson, and L. Shu, "Zero-shot out-of-distribution detection based on the pretrained model clip," in *AAAI*, 2022. [2](#), [5](#), [6](#)
- [40] H. Wang, Y. Li, H. Yao, and X. Li, "Clipn for zero-shot ood detection: Teaching clip to say no," in *ICCV*, 2023. [2](#), [5](#), [6](#), [7](#)
- [41] Y. Ming and Y. Li, "How does fine-tuning impact out-of-distribution detection for vision-language models?," *IJCV*, vol. 132, no. 2, pp. 596–609, 2024. [2](#), [6](#), [7](#), [9](#)
- [42] W. Tu, W. Deng, and T. Gedeon, "A closer look at the robustness of contrastive language-image pre-training (clip)," in *NeurIPS*, 2023. [2](#)
- [43] S. Park, J. Mok, D. Jung, S. Lee, and S. Yoon, "On the powerfulness of textual outlier exposure for visual ood detection," in *NeurIPS*, 2023. [2](#)
- [44] J. Nie, Y. Zhang, Z. Fang, T. Liu, B. Han, and X. Tian, "Out-of-distribution detection with negative prompts," in *ICLR*, 2023. [2](#), [6](#), [7](#), [10](#), [14](#)
- [45] X. Jiang, F. Liu, Z. Fang, H. Chen, T. Liu, F. Zheng, and B. Han, "Negative label guided ood detection with pretrained vision-language models," in *ICLR*, 2024. [2](#), [5](#), [6](#)
- [46] Y. Bai, Z. Han, C. Zhang, B. Cao, X. Jiang, and Q. Hu, "Id-like prompt learning for few-shot out-of-distribution detection," in *CVPR*, 2024. [2](#), [5](#), [6](#), [7](#)
- [47] T. Li, G. Pang, X. Bai, W. Miao, and J. Zheng, "Learning transferable negative prompts for out-of-distribution detection," in *CVPR*, 2024. [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [10](#), [13](#), [14](#)
- [48] C. Cao, Z. Zhong, Z. Zhou, Y. Liu, T. Liu, and B. Han, "Envisioning outlier exposure by large language models for out-of-distribution detection," in *ICML*, 2024. [2](#), [5](#), [6](#)
- [49] J. Liang, L. Sheng, Z. Wang, R. He, and T. Tan, "Realistic unsupervised CLIP fine-tuning with universal entropy optimization," in *ICML*, 2024. [2](#), [4](#), [8](#)
- [50] Y. Cao, X. Xu, J. Zhang, Y. Cheng, X. Huang, G. Pang, and W. Shen, "A survey on visual anomaly detection: Challenge, approach, and prospect," *arXiv preprint arXiv:2401.16402*, 2024. [2](#)
- [51] J. Liu, G. Xie, J. Wang, S. Li, C. Wang, F. Zheng, and Y. Jin, "Deep industrial image anomaly detection: A survey," *Machine Intelligence Research*, vol. 21, no. 1, pp. 104–135, 2024. [2](#)
- [52] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger, "Mvtec ad—a comprehensive real-world dataset for unsupervised anomaly detection," in *CVPR*, 2019. [3](#), [4](#), [5](#), [12](#)
- [53] Y. Zou, J. Jeong, L. Pemula, D. Zhang, and O. Dabeer, "Spot-the-difference self-supervised pre-training for anomaly detection and segmentation," in *ECCV*, 2022. [3](#), [4](#), [5](#), [12](#)
- [54] A. Palechor, A. Bhoomik, and M. Günther, "Large-scale open-set classification protocols for imagenet," in *WACV*, 2023. [3](#), [4](#), [5](#), [7](#), [14](#)
- [55] M. Salehi, H. Mirzaei, D. Hendrycks, Y. Li, M. H. Rohban, and M. Sabokrou, "A unified survey on anomaly, novelty, open-set, and out-of-distribution detection: Solutions and future challenges," *Transactions on Machine Learning Research*, 2022. [3](#), [4](#)
- [56] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, C. Li, J. Yang, H. Su, J. Zhu, *et al.*, "Grounding dino: Marrying dino with grounded pre-training for open-set object detection," *arXiv preprint arXiv:2303.05499*, 2023. [3](#), [10](#)
- [57] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, *et al.*, "Segment anything," in *ICCV*, 2023. [3](#), [10](#), [14](#)
- [58] H. Du, S. Zhang, B. Xie, G. Nan, J. Zhang, J. Xu, H. Liu, S. Leng, J. Liu, H. Fan, *et al.*, "Uncovering what why and how: A comprehensive benchmark for causation understanding of video anomaly," in *CVPR*, 2024. [4](#)
- [59] G. Zara, S. Roy, P. Rota, and E. Ricci, "Autolabel: Clip-based framework for open-set video domain adaptation," in *CVPR*, 2023. [4](#)
- [60] P. Wu, X. Zhou, G. Pang, Y. Sun, J. Liu, P. Wang, and Y. Zhang, "Open-vocabulary video anomaly detection," in *CVPR*, 2024. [4](#)
- [61] H. Deng, Z. Zhang, J. Bao, and X. Li, "Anovl: Adapting vision-language models for unified zero-shot anomaly localization," *arXiv preprint arXiv:2308.15939*, 2023. [4](#), [5](#), [6](#), [8](#)
- [62] X. Chen, J. Zhang, G. Tian, H. He, W. Zhang, Y. Wang, C. Wang, Y. Wu, and Y. Liu, "Clip-ad: A language-guided staged dual-path model for zero-shot anomaly detection," *arXiv preprint arXiv:2311.00453*, 2023. [4](#), [6](#), [8](#), [10](#)
- [63] X. Chen, Y. Han, and J. Zhang, "A zero-/few-shot anomaly classification and segmentation method for cvpr 2023 vand workshop challenge tracks 1&2: 1st place on zero-shot ad and 4th place on few-shot ad," *arXiv preprint arXiv:2305.17382*, 2023. [4](#), [5](#), [6](#), [8](#), [9](#), [13](#)
- [64] M. Tamura, "Random word data augmentation with clip for zero-shot anomaly detection," in *BMVC*, 2023. [4](#), [6](#), [8](#), [9](#)
- [65] Z. Gu, B. Zhu, G. Zhu, Y. Chen, H. Li, M. Tang, and J. Wang, "Filo: Zero-shot anomaly detection by fine-grained description and

- high-quality localization," *arXiv preprint arXiv:2404.13671*, 2024. 4, 5, 6, 8, 9, 13
- [66] Y. Li, A. Goodge, F. Liu, and C.-S. Foo, "Promptad: Zero-shot anomaly detection using text prompts," in *WACV*, 2024. 4
- [67] L. Deng, "The mnist database of handwritten digit images for machine learning research [best of the web]," *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012. 4
- [68] A. Krizhevsky, G. Hinton, *et al.*, "Learning multiple layers of features from tiny images," 2009. 4
- [69] A. Krizhevsky, V. Nair, and G. Hinton, "Cifar-10 and cifar-100 datasets," URL: <https://www.cs.toronto.edu/kriz/cifar.html>, vol. 6, no. 1, p. 1, 2009. 4
- [70] A. Torralba, R. Fergus, and W. T. Freeman, "80 million tiny images: A large data set for nonparametric object and scene recognition," *IEEE TPAMI*, vol. 30, no. 11, pp. 1958–1970, 2008. 4
- [71] Y. Shu, X. Guo, J. Wu, X. Wang, J. Wang, and M. Long, "Clipood: Generalizing clip to out-of-distributions," in *ICML*, 2023. 4, 13
- [72] K. Zhou, Z. Liu, Y. Qiao, T. Xiang, and C. C. Loy, "Domain generalization: A survey," *IEEE TPAMI*, vol. 45, no. 4, pp. 4396–4415, 2022. 4
- [73] Q. Yu, D. Ikami, G. Irie, and K. Aizawa, "Multi-task curriculum framework for open-set semi-supervised learning," in *ECCV*, 2020. 4
- [74] K. Saito, D. Kim, and K. Saenko, "Openmatch: Open-set semi-supervised learning with open-set consistency regularization," in *NeurIPS*, 2021. 4
- [75] K. Cao, M. Brbic, and J. Leskovec, "Open-world semi-supervised learning," in *ICLR*, 2022. 4
- [76] Y. Wang, W. Liu, X. Ma, J. Bailey, H. Zha, L. Song, and S.-T. Xia, "Iterative learning with open-set noisy labels," in *CVPR*, 2018. 4
- [77] K. Han, A. Vedaldi, and A. Zisserman, "Learning to discover novel visual categories via deep transfer clustering," in *CVPR*, 2019. 4
- [78] B. Zhao and K. Han, "Novel visual category discovery with dual ranking statistics and mutual knowledge distillation," in *NeurIPS*, 2021. 4
- [79] X. Jia, K. Han, Y. Zhu, and B. Green, "Joint representation learning and novel category discovery on single-and multi-modal data," in *ICCV*, 2021. 4
- [80] S. Vaze, K. Han, A. Vedaldi, and A. Zisserman, "Generalized category discovery," in *CVPR*, 2022. 4
- [81] K. Joseph, S. Paul, G. Aggarwal, S. Biswas, P. Rai, K. Han, and V. N. Balasubramanian, "Novel class discovery without forgetting," in *ECCV*, 2022. 4
- [82] G. Van Horn, O. Mac Aodha, Y. Song, Y. Cui, C. Sun, A. Shepard, H. Adam, P. Perona, and S. Belongie, "The inaturalist species classification and detection dataset," in *CVPR*, 2018. 5
- [83] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, "Places: A 10 million image database for scene recognition," *TPAMI*, vol. 40, no. 6, pp. 1452–1464, 2017. 5
- [84] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba, "Sun database: Large-scale scene recognition from abbey to zoo," in *CVPR*, 2010. 5
- [85] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi, "Describing textures in the wild," in *CVPR*, 2014. 5
- [86] M. C. Jung, H. Zhao, J. Dipnall, B. Gabbe, and L. Du, "Enhancing near ood detection in prompt learning: Maximum gains, minimal costs," *arXiv preprint arXiv:2405.16091*, 2024. 5, 14
- [87] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger, "Uninformed students: Student-teacher anomaly detection with discriminative latent embeddings," in *CVPR*, 2020. 5
- [88] T. Defard, A. Setkov, A. Loesch, and R. Audigier, "Padim: a patch distribution modeling framework for anomaly detection and localization," in *ICPR*, 2021. 5
- [89] C.-L. Li, K. Sohn, J. Yoon, and T. Pfister, "Cutpaste: Self-supervised learning for anomaly detection and localization," in *CVPR*, 2021. 5
- [90] P. Liznerski, L. Ruff, R. A. Vandermeulen, B. J. Franks, M. Kloft, and K.-R. Müller, "Explainable deep one-class classification," in *ICLR*, 2021. 5
- [91] J. Yi and S. Yoon, "Patch svdd: Patch-level svdd for anomaly detection and segmentation," in *ACCV*, 2020. 5
- [92] V. Zavrtnik, M. Kristan, and D. Skočaj, "Draem-a discriminatively trained reconstruction embedding for surface anomaly detection," in *ICCV*, 2021. 5
- [93] S. Thulasidasan, S. Thapa, S. Dhaubhadel, G. Chennupati, T. Bhattacharya, and J. Bilmes, "An effective baseline for robustness to distributional shift," *arXiv preprint arXiv:2105.07107*, 2021. 5
- [94] W.-H. Chu and K. M. Kitani, "Neural batch sampling with reinforcement learning for semi-supervised anomaly detection," in *ECCV*, 2020. 5
- [95] D. J. Atha and M. R. Jahanshahi, "Evaluation of deep learning approaches based on convolutional neural networks for corrosion detection," *Structural Health Monitoring*, 2018. 5
- [96] A. Miyai, Q. Yu, G. Irie, and K. Aizawa, "Zero-shot in-distribution detection in multi-object settings using vision-language foundation models," *arXiv preprint arXiv:2304.04521*, 2023. 6, 14
- [97] Y. Li, B. Xiong, G. Chen, and Y. Chen, "Setar: Out-of-distribution detection with selective low-rank approximation," *arXiv preprint arXiv:2406.12629*, 2024. 6
- [98] M. Lafon, E. Ramzi, C. Rambour, N. Audebert, and N. Thome, "Gallop: Learning global and local prompts for vision-language models," in *ECCV*, 2024. 6, 7, 14
- [99] X. Chen, Y. Li, and H. Chen, "Dual-adapter: Training-free dual adaptation for few-shot out-of-distribution detection," *arXiv preprint arXiv:2405.16146*, 2024. 6, 7, 9, 14
- [100] C. Ding and G. Pang, "Zero-shot out-of-distribution detection with outlier label exposure," in *IJCNN*, 2024. 6
- [101] D. Hendrycks and K. Gimpel, "A baseline for detecting misclassified and out-of-distribution examples in neural networks," in *ICLR*, 2017. 6, 14, 15
- [102] S. Liang, Y. Li, and R. Srikant, "Enhancing the reliability of out-of-distribution image detection in neural networks," in *ICLR*, 2018. 6, 14
- [103] K. Lee, K. Lee, H. Lee, and J. Shin, "A simple unified framework for detecting out-of-distribution samples and adversarial attacks," in *NeurIPS*, 2018. 6, 14
- [104] W. Liu, X. Wang, J. D. Owens, and Y. Li, "Energy-based out-of-distribution detection," in *NeurIPS*, 2020. 6, 14, 15
- [105] C. S. Sastry and S. Oore, "Detecting out-of-distribution examples with gram matrices," in *ICML*, 2020. 6, 14
- [106] H. Wang, W. Liu, A. Bocchieri, and Y. Li, "Can multi-label classification networks know what they don't know?," in *NeurIPS*, 2021. 6
- [107] J. Zhang, Q. Fu, X. Chen, L. Du, Z. Li, G. Wang, S. Han, D. Zhang, *et al.*, "Out-of-distribution detection based on in-distribution data patterns memorization with modern hopfield energy," in *ICLR*, 2023. 6
- [108] Y. Sun and Y. Li, "Dice: Leveraging sparsification for out-of-distribution detection," in *ECCV*, 2022. 6, 14, 15
- [109] Y. Sun, Y. Ming, X. Zhu, and Y. Li, "Out-of-distribution detection with deep nearest neighbors," in *ICML*, 2022. 6, 14, 15
- [110] Z. Lin, S. D. Roy, and Y. Li, "Mood: Multi-level out-of-distribution detection," in *CVPR*, 2021. 6
- [111] C. S. Sastry and S. Oore, "Detecting out-of-distribution examples with in-distribution examples and gram matrices," in *NeurIPS-W*, 2019. 6, 14
- [112] P. Sharma, N. Ding, S. Goodman, and R. Soicrut, "Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning," in *ACL*, 2018. 7
- [113] R. Zhang, W. Zhang, R. Fang, P. Gao, K. Li, J. Dai, Y. Qiao, and H. Li, "Tip-adapter: Training-free adaption of clip for few-shot classification," in *ECCV*, 2022. 7
- [114] Y. Zhang, W. Zhu, C. He, and L. Zhang, "Lapt: Label-driven automated prompt tuning for ood detection with vision-language models," in *ECCV*, 2024. 7
- [115] F. Lu, K. Zhu, K. Zheng, W. Zhai, and Y. Cao, "Likelihood-aware semantic alignment for full-spectrum out-of-distribution detection," *arXiv preprint arXiv:2312.01732*, 2023. 8, 13, 14
- [116] J. Yang, K. Zhou, and Z. Liu, "Full-spectrum out-of-distribution detection," *IJCV*, vol. 131, no. 10, pp. 2607–2622, 2023. 8, 14
- [117] D. Hendrycks and T. Dietterich, "Benchmarking neural network robustness to common corruptions and perturbations," in *ICLR*, 2019. 8
- [118] D. Hendrycks, S. Basart, N. Mu, S. Kadavath, F. Wang, E. Dorundo, R. Desai, T. Zhu, S. Parajuli, M. Guo, *et al.*, "The many faces of robustness: A critical analysis of out-of-distribution generalization," in *ICCV*, 2021. 8
- [119] J. Zhang, J. Yang, P. Wang, H. Wang, Y. Lin, H. Zhang, Y. Sun, X. Du, K. Zhou, W. Zhang, Y. Li, Z. Liu, Y. Chen, and H. Li, "Openood v1.5: Enhanced benchmark for out-of-distribution detection," *arXiv preprint arXiv:2306.09301*, 2023. 8, 15

- [120] B. Recht, R. Roelofs, L. Schmidt, and V. Shankar, "Do imagenet classifiers generalize to imagenet?," in *ICML*, 2019. 8
- [121] T. Huang, J. Chu, and F. Wei, "Unsupervised prompt learning for vision-language models," *arXiv preprint arXiv:2204.03649*, 2022. 8
- [122] M. Shu, W. Nie, D.-A. Huang, Z. Yu, T. Goldstein, A. Anandkumar, and C. Xiao, "Test-time prompt tuning for zero-shot generalization in vision-language models," in *NeurIPS*, 2022. 8
- [123] K. Tanwisuth, S. Zhang, H. Zheng, P. He, and M. Zhou, "Pouf: Prompt-oriented unsupervised fine-tuning for large pre-trained models," in *ICML*, 2023. 8
- [124] Z. Zhou, M. Yang, J.-X. Shi, L.-Z. Guo, and Y.-F. Li, "Decoop: Robust prompt tuning with out-of-distribution detection," in *ICML*, 2024. 8
- [125] P. Isola, J. J. Lim, and E. H. Adelson, "Discovering states and transformations in image collections," in *CVPR*, 2015. 8
- [126] N. Belton, M. T. Hagos, A. Lawlor, and K. M. Curran, "Fewsome: One-class few shot anomaly detection with siamese networks," in *CVPR*, 2023. 9
- [127] C. Huang, H. Guan, A. Jiang, Y. Zhang, M. Spratling, and Y.-F. Wang, "Registration based few-shot anomaly detection," in *ECCV*, 2022. 9
- [128] J. Liao, X. Xu, M. C. Nguyen, A. Goodge, and C. S. Foo, "Coft-ad: Contrastive fine-tuning for few-shot anomaly detection," *IEEE TIP*, vol. 33, pp. 2090–2103, 2024. 9
- [129] E. Schwartz, A. Arbelle, L. Karlinsky, S. Harary, F. Scheidegger, S. Doveh, and R. Giryes, "Maeday: Mae for few-and zero-shot anomaly-detection," *Computer Vision and Image Understanding*, p. 103958, 2024. 9
- [130] S. Sheynin, S. Benaïm, and L. Wolf, "A hierarchical transformation-discriminating generative model for few shot anomaly detection," in *ICCV*, 2021. 9
- [131] Z. Wang, Y. Zhou, R. Wang, T.-Y. Lin, A. Shah, and S. N. Lim, "Few-shot fast-adaptive anomaly detection," in *NeurIPS*, 2022. 9
- [132] J.-C. Wu, D.-J. Chen, C.-S. Fuh, and T.-L. Liu, "Learning unsupervised metaformer for anomaly detection," in *ICCV*, 2021. 9
- [133] G. Xie, J. Wang, J. Liu, F. Zheng, and Y. Jin, "Pushing the limits of fewshot anomaly detection in industry vision: Graphcore," in *ICLR*, 2023. 9
- [134] K. Roth, L. Pemula, J. Zepeda, B. Schölkopf, T. Brox, and P. Gehler, "Towards total recall in industrial anomaly detection," in *CVPR*, 2022. 9
- [135] Y. Cao, X. Xu, C. Sun, Y. Cheng, Z. Du, L. Gao, and W. Shen, "Segment any anomaly without training via hybrid prompt regularization," *arXiv preprint arXiv:2305.10724*, 2023. 10
- [136] C. Li, L. Qi, and X. Geng, "A sam-guided two-stream lightweight model for anomaly detection," *arXiv preprint arXiv:2402.19145*, 2024. 10
- [137] L. Hua, Y. Luo, Q. Qi, and J. Long, "Medicalclip: Anomaly-detection domain generalization with asymmetric constraints," *Biomolecules*, vol. 14, no. 5, p. 590, 2024. 10
- [138] X. Zhang, M. Xu, D. Qiu, R. Yan, N. Lang, and X. Zhou, "Mediclip: Adapting clip for few-shot medical image anomaly detection," *arXiv preprint arXiv:2405.11315*, 2024. 10
- [139] X. Du, Z. Wang, M. Cai, and Y. Li, "Vos: Learning what you don't know by virtual outlier synthesis," in *ICLR*, 2022. 10
- [140] X. Du, X. Wang, G. Gozum, and Y. Li, "Unknown-aware object detection: Learning what you don't know from videos in the wild," in *CVPR*, 2022. 10
- [141] X. Du, G. Gozum, Y. Ming, and Y. Li, "Siren: Shaping representations for detecting out-of-distribution objects," in *NeurIPS*, 2022. 10
- [142] Y. Cao, X. Xu, C. Sun, X. Huang, and W. Shen, "Towards generic anomaly detection and understanding: Large-scale visual-linguistic model (gpt-4v) takes the lead," *arXiv preprint arXiv:2311.02782*, 2023. 11, 12
- [143] A. Miyai, J. Yang, J. Zhang, Y. Ming, Q. Yu, G. Irie, Y. Li, H. Li, Z. Liu, and K. Aizawa, "Unsolvable problem detection: Evaluating trustworthiness of vision language models," *arXiv preprint arXiv:2403.20331*, 2024. 11, 12
- [144] Z. Gu, B. Zhu, G. Zhu, Y. Chen, M. Tang, and J. Wang, "Anomalygpt: Detecting industrial anomalies using large vision-language models," in *AAAI*, 2024. 11, 12
- [145] Y. Li, H. Wang, S. Yuan, M. Liu, D. Zhao, Y. Guo, C. Xu, G. Shi, and W. Zuo, "Myriad: Large multimodal model by applying vision experts for industrial anomaly detection," *arXiv preprint arXiv:2310.19070*, 2023. 11
- [146] L. Chen, S. Li, J. Yan, H. Wang, K. Gunaratna, V. Yadav, Z. Tang, et al., "Alpagasus: Training a better alpaca with fewer data," *arXiv preprint arXiv:2307.08701*, 2023. 11
- [147] W.-L. Chiang, Z. Li, Z. Lin, Y. Sheng, et al., "Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality." <https://lmsys.org/blog/2023-03-30-vicuna>, 2023. Accessed: 2024-07-02. 11
- [148] Z. Li, P. Xu, F. Liu, and H. Song, "Towards understanding in-context learning with contrastive demonstrations and saliency maps," *arXiv preprint arXiv:2307.05052*, 2023. 11
- [149] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al., "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023. 11
- [150] J. Wei, J. Wei, Y. Tay, D. Tran, A. Webson, Y. Lu, X. Chen, H. Liu, D. Huang, D. Zhou, et al., "Larger language models do in-context learning differently," *arXiv preprint arXiv:2303.03846*, 2023. 11
- [151] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, et al., "A survey of large language models," *arXiv preprint arXiv:2303.18223*, 2023. 11
- [152] A. Awadalla, I. Gao, J. Gardner, J. Hessel, Y. Hanafy, W. Zhu, K. Marathe, Y. Bitton, S. Gadre, S. Sagawa, et al., "Openflamingo: An open-source framework for training large autoregressive vision-language models," *arXiv preprint arXiv:2308.01390*, 2023. 11
- [153] S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T. Lee, Y. Li, S. Lundberg, et al., "Sparks of artificial general intelligence: Early experiments with gpt-4," *arXiv preprint arXiv:2303.12712*, 2023. 11
- [154] W. Dai, J. Li, D. Li, A. Tiong, J. Zhao, W. Wang, B. Li, P. Fung, and S. Hoi, "InstructBLIP: Towards general-purpose vision-language models with instruction tuning," in *NeurIPS*, 2023. 11
- [155] W. Wang, Q. Lv, W. Yu, W. Hong, J. Qi, Y. Wang, J. Ji, Z. Yang, L. Zhao, X. Song, et al., "Cogvlm: Visual expert for pretrained language models," *arXiv preprint arXiv:2311.03079*, 2023. 11
- [156] Q. Ye, H. Xu, G. Xu, J. Ye, M. Yan, Y. Zhou, J. Wang, A. Hu, P. Shi, Y. Shi, et al., "mplug-owl: Modularization empowers large language models with multimodality," *arXiv preprint arXiv:2304.14178*, 2023. 11
- [157] D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny, "Minigtpt-4: Enhancing vision-language understanding with advanced large language models," in *ICLR*, 2024. 11
- [158] B. Li, Y. Zhang, L. Chen, J. Wang, J. Yang, and Z. Liu, "Otter: A multi-modal model with in-context instruction tuning," *arXiv preprint arXiv:2305.03726*, 2023. 11
- [159] J. Lin, H. Yin, W. Ping, Y. Lu, P. Molchanov, A. Tao, H. Mao, J. Kautz, M. Shoenybi, and S. Han, "Vila: On pre-training for visual language models," in *CVPR*, 2024. 11
- [160] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh, "Vqa: Visual question answering," in *ICCV*, 2015. 11
- [161] F. Liu, H. Tan, and C. Tensmeyer, "Documentclip: Linking figures and main body text in reflowed documents," in *ICPRAI*, 2024. 11
- [162] Y. Liu, H. Duan, Y. Zhang, B. Li, S. Zhang, W. Zhao, Y. Yuan, J. Wang, C. He, Z. Liu, et al., "Mmbench: Is your multi-modal model an all-around player?," in *ECCV*, 2024. 11
- [163] X. Yue, Y. Ni, K. Zhang, T. Zheng, R. Liu, G. Zhang, S. Stevens, D. Jiang, W. Ren, Y. Sun, et al., "Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi," in *CVPR*, 2024. 11, 15
- [164] F. Wang, X. Fu, J. Y. Huang, Z. Li, Q. Liu, X. Liu, M. D. Ma, N. Xu, W. Zhou, K. Zhang, et al., "Muirbench: A comprehensive benchmark for robust multi-image understanding," *arXiv preprint arXiv:2406.09411*, 2024. 11, 15
- [165] J. Liu, Y. Yuan, J. Hao, F. Ni, L. Fu, Y. Chen, and Y. Zheng, "Enhancing robotic manipulation with ai feedback from multimodal large language models," *arXiv preprint arXiv:2402.14245*, 2024. 11
- [166] Y. Li, W. Zhang, K. Chen, Y. Liu, P. Li, R. Gao, L. Hong, M. Tian, X. Zhao, Z. Li, et al., "Automated evaluation of large vision-language models on self-driving corner cases," *arXiv preprint arXiv:2404.10595*, 2024. 11
- [167] OpenAI, "Gpt-4v(ision) system card." https://cdn.openai.com/papers/GPTV_System_Card.pdf, 2023. 11
- [168] P. Bergmann, X. Jin, D. Sattlegger, and C. Steger, "The mvtec 3d-ad dataset for unsupervised 3d anomaly detection and localization," in *VISAPP*, 2022. 12

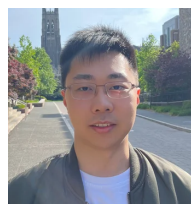
- [169] D. S. Kermany, M. Goldbaum, W. Cai, *et al.*, "Identifying medical diagnoses and treatable diseases by image-based deep learning," *Cell*, vol. 172, no. 5, pp. 1122–1131, 2018. **12**
- [170] K. Felipe, "Head ct - hemorrhage." <https://www.kaggle.com/datasets/felipekitamura/headct-hemorrhage>, 2018. **12**
- [171] P. Bergmann, K. Batzner, M. Fauser, D. Sattlegger, and C. Steger, "Beyond dents and scratches: Logical constraints in unsupervised anomaly detection and localization," *IJCV*, vol. 130, no. 4, pp. 947–969, 2022. **12**
- [172] W. Sultani, C. Chen, and M. Shah, "Real-world anomaly detection in surveillance videos," in *CVPR*, 2018. **12**
- [173] C. Kay, "Accident detection from cctv footage." <https://www.kaggle.com/datasets/ckay16/accident-detection-fromcctv-footage>, 2018. Kaggle dataset. **12**
- [174] S. E. User, "Simple outlier detection for time series." <https://stats.stackexchange.com/questions/427327/simple-outlier-detection-for-time-series>, 2021. Accessed: 2024-07-02. **12**
- [175] M. Xu, Z. Zhang, F. Wei, H. Hu, and X. Bai, "Side adapter network for open-vocabulary semantic segmentation," in *CVPR*, 2023. **13**
- [176] X. Gu, T.-Y. Lin, W. Kuo, and Y. Cui, "Open-vocabulary object detection via vision and language knowledge distillation," in *ICLR*, 2022. **13**
- [177] Y. Shu, Z. Cao, C. Wang, J. Wang, and M. Long, "Open domain generalization with domain-augmented meta-learning," in *CVPR*, 2021. **13**
- [178] K. Katsumata, I. Kishida, A. Amma, and H. Nakayama, "Open-set domain generalization via metric learning," in *ICIP*, 2021. **13**
- [179] R. Zhu and S. Li, "Crossmatch: Cross-classifier consistency regularization for open-set single domain generalization," in *ICLR*, 2021. **13**
- [180] M. Noguchi and S. Shirakawa, "Simple domain generalization methods are strong baselines for open domain generalization," *arXiv preprint arXiv:2303.18031*, 2023. **13**
- [181] Z. Chen, W. Wang, Z. Zhao, F. Su, A. Men, and H. Meng, "Practicaldgl: Perturbation distillation on vision-language models for hybrid domain generalization," in *CVPR*, 2024. **13**
- [182] M. Singha, A. Jha, S. Bose, A. Nair, M. Abdar, and B. Banerjee, "Unknown prompt the only lacuna: Unveiling clip's potential for open domain generalization," in *CVPR*, 2024. **13**
- [183] Y. Sun, C. Guo, and Y. Li, "React: Out-of-distribution detection with rectified activations," in *NeurIPS*, 2021. **14**
- [184] X. Dong, J. Guo, A. Li, W.-T. Ting, C. Liu, and H. Kung, "Neural mean discrepancy for efficient out-of-distribution detection," in *CVPR*, 2022. **14**
- [185] Y. Ming, H. Yin, and Y. Li, "On the impact of spurious correlation for out-of-distribution detection," in *AAAI*, 2022. **14**
- [186] G. Chen, W. Yao, X. Song, X. Li, Y. Rao, and K. Zhang, "Plot: Prompt learning with optimal transport for vision-language models," in *ICLR*, 2023. **14**
- [187] A. Bulat and G. Tzimiropoulos, "Lasp: Text-to-text optimization for language-aware soft prompting of vision & language models," in *CVPR*, 2023. **14**
- [188] Y. Lu, J. Liu, Y. Zhang, Y. Liu, and X. Tian, "Prompt distribution learning," in *CVPR*, 2022. **14**
- [189] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. De Larousilhe, A. Gesmundo, M. Attariyan, and S. Gelly, "Parameter-efficient transfer learning for nlp," in *ICML*, 2019. **14**
- [190] M. Jia, L. Tang, B.-C. Chen, C. Cardie, S. Belongie, B. Hariharan, and S.-N. Lim, "Visual prompt tuning," in *ECCV*, 2022. **14**
- [191] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," in *ICLR*, 2022. **14, 15**
- [192] V. Udandarao, A. Gupta, and S. Albanie, "Sus-x: Training-free name-only transfer of vision-language models," in *ICCV*, 2023. **14**
- [193] Y. Ming and Y. Li, "Understanding retrieval-augmented task adaptation for vision-language models," in *ICML*, 2024. **14**
- [194] H. Bai, G. Canal, X. Du, J. Kwon, R. D. Nowak, and Y. Li, "Feed two birds with one scone: Exploiting wild data for both out-of-distribution generalization and detection," in *ICML*, 2023. **14**
- [195] M. U. Khattak, S. T. Wasim, M. Naseer, S. Khan, M.-H. Yang, and F. S. Khan, "Self-regulating prompts: Foundational model adaptation without forgetting," in *ICCV*, 2023. **14**
- [196] C. Leys, O. Klein, Y. Dominicy, and C. Ley, "Detecting multivariate outliers: Use a robust variant of the mahalanobis distance," *Journal of Experimental Social Psychology*, vol. 74, pp. 150–156, 2018. **15**
- [197] H. Wang, Z. Li, L. Feng, and W. Zhang, "Vim: Out-of-distribution with virtual-logit matching," in *CVPR*, 2022. **15**
- [198] A. Miyai, Q. Yu, G. Irie, and K. Aizawa, "Can pre-trained networks detect familiar out-of-distribution data?," *arXiv preprint arXiv:2310.00847*, 2023. **15**
- [199] X. Chen, H. Fan, R. Girshick, and K. He, "Improved baselines with momentum contrastive learning," *arXiv preprint arXiv:2003.04297*, 2020. **15**
- [200] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, "Emerging properties in self-supervised vision transformers," in *ICCV*, 2021. **15**
- [201] J. Dong, Y. Gao, H. Zhou, J. Cen, Y. Yao, S. Yoon, and P. D. Sun, "Towards few-shot out-of-distribution detection," *arXiv preprint arXiv:2311.12076*, 2023. **15**
- [202] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," in *ICLR*, 2021. **15**
- [203] Y. Zhang, K. Zhou, and Z. Liu, "Neural prompt search," *arXiv preprint arXiv:2206.04673*, 2022. **15**
- [204] E. Baek, K. Park, J. Kim, and H.-S. Kim, "Unexplored faces of robustness and out-of-distribution: Covariate shifts in environment and sensor domains," in *CVPR*, 2024. **15**
- [205] P. W. Koh, S. Sagawa, H. Marklund, S. M. Xie, M. Zhang, A. Balsubramani, W. Hu, M. Yasunaga, R. L. Phillips, I. Gao, *et al.*, "Wilds: A benchmark of in-the-wild distribution shifts," in *ICML*, 2021. **15**
- [206] L. Cultrera, L. Seidenari, and A. Del Bimbo, "Leveraging visual attention for out-of-distribution detection," in *ICCV*, 2023. **15**
- [207] Z. Hong, Y. Yue, Y. Chen, H. Lin, Y. Luo, *et al.*, "Out-of-distribution detection in medical image analysis: A survey," *arXiv preprint arXiv:2404.18279*, 2024. **15**
- [208] P. Lu, H. Bansal, T. Xia, J. Liu, C. Li, H. Hajishirzi, H. Cheng, K.-W. Chang, M. Galley, and J. Gao, "Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts," in *ICLR*, 2024. **15**



Atsuyuki Miyai (Student Member, IEEE) received his B.E. degree in Information and Communication Engineering in 2022 and his M.E. degree in Interdisciplinary Studies in Information Science in 2024 from The University of Tokyo. He is currently a Ph.D. student with the Department of Information and Communication Engineering at The University of Tokyo. His research topic is the safety of vision language models and foundation models. He is an IEEE Student Branch Chair at The University of Tokyo.



Jingyang Yang is currently a PhD student at the College of Computing and Data Science (CCDS) at Nanyang Technological University (NTU), Singapore, working in the MMLab@NTU under the supervision of Dr. Ziwei Liu. His research interests include visual generalist models and AI safety for foundation models. He has published over 20 papers in relevant fields at top-tier conferences and journals, including CVPR, ICCV, ECCV, NeurIPS, ICLR, and IJCV. He serves as a reviewer of CVPR, ICCV, ECCV, NeurIPS, ICLR, etc., and is an outstanding reviewer in ICCV 2021 and CVPR 2024.



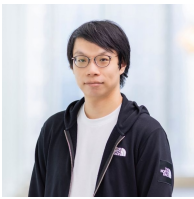
Jingyang Zhang received his B.E. degree in Electronic Engineering from Tsinghua University, Beijing, China, in 2019. Since then, he has been pursuing the Ph.D. degree at Dept. of Electrical and Computer Engineering at Duke University, supervised by Dr. Yiran Chen and Dr. Hai (Helen) Li. His research spans from robustness of deep learning-based vision systems to (more recently) generative AI and multi-modal LLMs.



Yifei Ming received his Ph.D. in Computer Science from the University of Wisconsin-Madison in 2024, advised by Dr. Yixuan Li. He is currently a Research Scientist at Salesforce AI Research. His research interests encompass reliable machine learning, vision-language models, large language models, and responsible foundation models. He has a prolific publication record and serves on the program committees of leading machine learning conferences and journals, including NeurIPS, ICLR, ICML, EMNLP, and IJCV.



Yueqian Lin received a B.S. degree in Data Science from Duke Kunshan University in 2024, graduating summa cum laude with signature work distinction. He is currently pursuing a Ph.D. in Electrical and Computer Engineering at Duke University under the supervision of Dr. Yiran Chen and Dr. Hai (Helen) Li. His research focuses on robustness in generative models.



Qing Yu (Member, IEEE) received his B.E. degree in Information and Communication Engineering and the M.E. degree in Interdisciplinary Studies in Information Science from The University of Tokyo in 2018 and 2020, respectively. He completed his Ph.D. in Information Science and Technology from The University of Tokyo in 2023. He is currently a Research Scientist at LY Corporation, Japan. His research interests include multimodal recognition and generation.



Go Irie (Member, IEEE) received his B.E. and M.E. in System Engineering from Keio University, Japan, in 2004 and 2006, respectively. He completed his Ph.D. in Information Science and Technology from The University of Tokyo in 2011. He is currently an Associate Professor at the Department of Information and Computer Technology, Tokyo University of Science, Japan. He was a research scientist at NTT Corporation, Japan, from 2006 to 2022, and a Visiting Research Scholar at Columbia University from 2012

to 2013. His research interests include pattern recognition, machine learning and media understanding.



Shafiq Joty is currently a Research Director at Salesforce Research (Palo Alto, USA), where he oversees the NLP group's efforts in large language modeling (LLM) and generative AI. He also holds the position of a tenured Associate Professor (currently on leave) in the School of Computer Science and Engineering (SCSE) at NTU, Singapore. He was a founding manager of the Salesforce Research Asia (Singapore) lab. His research has contributed to over 30+ patents and 145+ papers in top-tier NLP and

ML conferences and journals. He has served as the Program Chair of SIGDIAL-2023, as a member of the best paper award committees for ICLR-23 and NAACL-22, and in the capacity of a (senior) area chair for many of the leading NLP and ML conferences.



Yixuan Li (Member, IEEE) received the PhD degree from Cornell University, in 2017, advised by John E. Hopcroft. She is an assistant professor in the Department of Computer Sciences, the University of Wisconsin-Madison. Subsequently, she was a postdoctoral scholar in the computer science department, Stanford University. Her research focuses on the algorithmic and theoretical foundations of learning in open worlds. She has served as Area Chair for ICLR, NeurIPS, ICML, and Program Chair for Workshop on Uncertainty

and Robustness in Deep Learning. She is the recipient of the AFOSR Young Investigator Program (YIP) award, NSF CAREER award, MIT Technology Review Innovator Under 35, Forbes 30 Under 30 in Science, and multiple faculty research awards from Google, Meta, and Amazon. Her works received a NeurIPS Outstanding Paper Award, and an ICLR Outstanding Paper Award Honorable Mention, in 2022.



Hai (Helen) Li (Fellow, IEEE) received the Ph.D. degree from Purdue University in 2004. Dr. Li is currently the Clare Boothe Luce Professor and Department Chair of the Electrical and Computer Engineering Department at Duke University. Her current research interests include neuromorphic circuits and systems for brain-inspired computing, machine learning acceleration and trustworthy AI, conventional and emerging memory design and architecture, and software and hardware co-design. Dr. Li is a recipient of the NSF Career Award (2012), DARPA Young Faculty Award (2013), TUM-IAS Hans Fischer Fellowship from Germany (2017), and ELATE Fellowship (2020). She received 9 best paper awards and additional 9 best paper nominations from international conferences. Dr. Li is a Distinguished Lecturer of the IEEE CAS Society (2018-2019) and a Distinguished Speaker of ACM (2017-2020).



Ziwei Liu (Member, IEEE) is currently a nanyang assistant professor with Nanyang Technological University, Singapore. His research interests include computer vision machine learning and computer graphics. He has published extensively on top-tier conferences and journals in relevant fields, including CVPR, ICCV, ECCV, NeurIPS, ICLR, ICML, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *ACM Transactions on Graphics and Nature - Machine Intelligence*. He is the recipient of ICCV Young Researcher

Award, HKSTP Best Paper Award, CVPR Best Paper Award Candidate, ICBS Frontiers of Science Award and MIT Technology Review Innovators under 35 Asia Pacific. He serves as an area chair of CVPR, ICCV, ECCV, NeurIPS and ICLR, as well as an associate editor of *International Journal of Computer Vision*.



Toshihiko Yamasaki (Member, IEEE) received the Ph.D. degree from The University of Tokyo. He is currently a Professor at the Department of Information and Communication Engineering, Graduate School of Information Science and Technology, The University of Tokyo. He was a JSPS Fellow for Research Abroad and a visiting scientist at Cornell University from 2011 to 2013. His current research interests are computer vision, multimedia, pattern recognition, machine learning. He is a pioneer of attractiveness computing.

puting.



Kiyoharu Aizawa (Fellow, IEEE) received the B.E., the M.E., and the Dr.Eng. degrees in Electrical Engineering all from the University of Tokyo, in 1983, 1985, 1988, respectively. He is a professor with the Department of Information and Communication Engineering, and director of VR center, University of Tokyo. He was a visiting assistant professor with the University of Illinois from 1990 to 1992. His research fields are multimedia, image processing, and computer vision, with a particular interest in interdisciplinary

and cross-disciplinary issues. He received the 1990, 1998 Best Paper Awards, the 1991 Achievement Award, 1999 Electronics Society Award from IEICE Japan, and the 1998 Fujio Frontier Award, the 2002 and 2009 Best Paper Award, and 2013 Achievement award from ITE Japan. He received the IBM Japan Science Prize in 2002. He is on Editorial Board of ACM TOMM. He served as the Editor-in-Chief of Journal of ITE Japan, an Associate Editor of IEEE TIP, TCSVT, TMM, and MultiMedia. He has served a number of international and domestic conferences; he was a General co-Chair of ACM Multimedia 2012 and ICMR2018. He is a Fellow of IEEE, IEICE, ITE and a member of Science Council of Japan.