

Synthesizing Near-Boundary OOD Samples for Out-of-Distribution Detection

Jinglun Li^{1,3}, Kaixun Jiang¹, Zhaoyu Chen¹, Bo Lin³, Yao Tang³, Weifeng Ge^{2†}, Wenqiang Zhang^{1,2†}

¹College of Intelligent Robotics and Advanced Manufacturing, Fudan University, Shanghai

²Shanghai Key Lab of Intelligent Information Processing,

College of Computer Science and Artificial Intelligence, Fudan University, Shanghai

³JIIOV Technology, Beijing

{jinglunli21, kxjiang22}@m.fudan.edu.cn, zhaoyuchen20@fudan.edu.cn,

{bo.lin, yao.tang}@jiiiov.com,

weifeng.ge.ic@gmail.com, wqzhang@fudan.edu.cn

Abstract

*Pre-trained vision-language models have exhibited remarkable abilities in detecting out-of-distribution (OOD) samples. However, some challenging OOD samples, which lie close to in-distribution (InD) data in image feature space, can still lead to misclassification. The emergence of foundation models like diffusion models and multimodal large language models (MLLMs) offers a potential solution to this issue. In this work, we propose **SynOOD**, a novel approach that harnesses foundation models to generate synthetic, challenging OOD data for fine-tuning CLIP models, thereby enhancing boundary-level discrimination between InD and OOD samples. Our method uses an iterative in-painting process guided by contextual prompts from MLLMs to produce nuanced, boundary-aligned OOD samples. These samples are refined through noise adjustments based on gradients from OOD scores like the energy score, effectively sampling from the InD/OOD boundary. With these carefully synthesized images, we fine-tune the CLIP image encoder and negative label features derived from the text encoder to strengthen connections between near-boundary OOD samples and a set of negative labels. Finally, SynOOD achieves state-of-the-art performance on the large-scale ImageNet benchmark, with minimal increases in parameters and runtime. Our approach significantly surpasses existing methods, improving AUROC by 2.80% and reducing FPR95 by 11.13%. Codes are available in <https://github.com/Jarvisgivemeasuit/SynOOD>.*

1. Introduction

Modern deep neural networks [8, 18, 33, 57] deployed in open-world scenarios inevitably encounter out-of-

distribution (OOD) samples, which can pose significant security risks. Accurate identification of OOD data is crucial to mitigate these threats. Traditional vision-based OOD detection methods [15, 21, 26, 28, 31, 46, 47] often rely solely on a single image domain. Recent research [22, 27, 35, 36, 54] in pre-trained visual-language models [25, 41] has demonstrated significant improvements in OOD detection by effectively employing both visual and language information. In particular some CLIP-based methods [22, 35, 54], such as NegLabel [22], enhance OOD detection by introducing potential OOD text labels, denoted as negative labels, that lie outside the in-distribution (InD) label space. However, a significant challenge remains in accurately identifying hard OOD samples near the InD/OOD boundary, as these samples often appear visually similar to InD instances, making them difficult to classify using CLIP-based methods directly. CLIP-based methods show that OOD samples situated near the InD/OOD boundary tend to align more closely with InD labels, as images are typically more densely packed in the feature space than labels, limiting the model’s ability to establish clear semantic alignment, this is illustrated in Fig. 1(a). Consequently, this mismatch reduces the reliability of CLIP [41] in detecting boundary OOD samples, particularly those closely resembling the InD distribution.

A promising approach to improve OOD detection is to effectively map ambiguous samples near the InD/OOD boundary to either InD or negative labels. However, fine-tuning CLIP models for this purpose has been challenging due to a lack of suitable data. Recent advancements in multimodal large language models (MLLMs) [1, 25, 29, 38] and diffusion models [40, 43] offer powerful generative capabilities, though their application to task-oriented data generation remains relatively unexplored. To fill this gap, we propose a novel iterative generative approach that utilizes

[†]indicates corresponding authors.

the contextual understanding of an MLLM and the sophisticated image synthesis capabilities of a diffusion model. This integration enables the creation of realistic, boundary-aligned OOD samples that are visually similar to InD data while remaining sufficiently distinct. By generating these nuanced near-boundary OOD samples, our approach provides the CLIP model with more challenging data for fine-tuning, achieving a more accurate separation between InD and OOD samples.

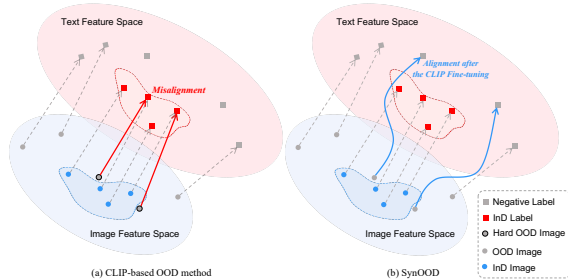


Figure 1. (a) illustrates a simplified example highlighting the limitations of CLIP-based OOD methods, where challenging OOD samples are misclassified due to CLIP models’ limited fine-grained discrimination. (b) Our proposed method addresses this limitation by generating challenging data to fine-tune the CLIP models, building strong connections between confusing OOD samples and their corresponding negative labels.

Our method begins by using a language model to extract all detected contextual elements within an InD image. For example, in an image labeled “panda,” the language model may detect contextual elements like “bamboo,” “tourist,” “leaf,” and “railing,” which are commonly associated with the primary subject but not central to it. These elements then serve as prompts for an in-painting diffusion model. Rather than relying on predefined masks, we employ an iterative generative process to guide the model in creating images that remain visually similar to the InD data while representing OOD content. In each iteration, the generated image is evaluated using an OOD detection model, and the resulting OOD score informs a gradient. This gradient is backpropagated through the diffusion model, updating the noise to iteratively adjust the image. Over time, the model gradually replaces primary subject elements, such as the “panda,” with background elements from the identified list. This controlled transformation subtly shifts the image’s focus away from the core theme, allowing it to appear distinct from InD samples without losing its underlying visual similarities. By iteratively integrating contextual elements as the main focus while maintaining the original style and setting, the resulting synthetic images closely resemble InD examples in appearance but remain distinctly OOD, aligning with the theoretical principles in [12, 63].

In this work, we propose **SynOOD**, a novel method

that iteratively generates near-boundary data to fine-tune the CLIP models for enhancing OOD detection performance. Specifically, our method contains three components: an iterative generative process, fine-tuning the CLIP image encoder with a projection layer, and refining negative label features derived from the CLIP text encoder. This process significantly boosts the model’s capacity to distinguish between InD and OOD samples, this is illustrated in Fig. 1(b). By integrating these processes, SynOOD offers a robust and effective approach to OOD detection. Our contributions are summarized as follows:

- We propose **SynOOD**, a novel framework for OOD detection that generates challenging, near-boundary OOD samples to the fine-tune CLIP models, enhancing to detect difficult OOD cases close to the InD/OOD boundary.
- We introduce a generation process that iterative synthesizes near-boundary OOD samples using foundation models, guided by OOD gradient information. This process yields high-quality, nuanced data that enhances CLIP to strengthen connections between challenging OOD samples and negative labels.
- Extensive experiments show that SynOOD achieves state-of-the-art performance on widely used large-scale benchmarks, with minimal increases in parameters and runtime, outperforming existing methods by improving AUROC by 2.80% and reducing FPR95 by 11.13%.

2. Related Work

OOD detection with visual modal. Single-modal visual OOD detection methods include: (1) Logit-based approaches, which compute OOD scores from network logits. MSP [15] uses the maximum logit, while ODIN [28] enhances separation via input perturbations and logit rescaling. ReAct [47] further reduces overconfidence by adjusting activation logits. (2) Distance-based methods, which use feature distances between InD and OOD samples as OOD scores. Gaussian discriminant analysis [24, 55] and metrics like cosine similarity [4, 59], Euclidean distance [19], and RBF kernels [50] are commonly employed. (3) Gradient-based methods, such as GradNorm [21], leverage classifier gradients to distinguish InD from OOD samples using gradient-based features.

OOD detection with multi-modal models. Leveraging textual information alongside visual data for OOD detection has become increasingly popular due to its strong performance. Fort et al.[13] pioneered this direction by using class names of potential outliers as input to image-text encoders like CLIP, improving OOD detection. MCM[35] is an effective post-hoc method that uses the maximum predicted softmax value from a vision-language model as the OOD score. More recently, CLIPN [54] proposed using a text encoder to identify OOD samples by comparing similarity discrepancies between two text encoders and a

frozen image encoder. Building on this, LSN [37] introduced negative classifiers with learned prompts to detect images outside specific categories. NPOS [49] generates synthetic OOD data to better define decision boundaries between InD and OOD samples. LAPT [62] automates prompt tuning for vision-language models, reducing manual effort. DreamOOD [10] learns a text-conditioned latent space to generate diverse OOD samples by decoding low-likelihood embeddings into images. NegLabel [22] selects potential OOD labels from semantically related WordNet [34] terms outside the InD label space, using a pre-trained vision-language model like CLIP to classify images as InD or OOD.

3. Method

3.1. OOD Detection Setup

Given a training set $\mathcal{D}^{\text{in}} = \{(x_i, y_i)\}_{i=1}^n$, where $x_i \in \mathbb{R}^{3 \times H \times W}$ is the 3-channel image of size $H \times W$, $y_i \in \mathcal{Y}$ denotes one of C InD categories, and n is the number of samples, our target is to develop an OOD detector $G(x)$ solely based on \mathcal{D}^{in} . When applied to a test image set $\mathcal{X} = \{x_i\}_{i=1}^K$, the detector $G(x)$ should output a binary classification using a score function $S(x)$:

$$G(x) = \begin{cases} \text{InD}, & \text{if } S(x) \geq \eta; \\ \text{OOD}, & \text{if } S(x) < \eta, \end{cases} \quad (1)$$

where η is a threshold parameter. We follow Jiang et al. [22] and employ a negative label set $\mathcal{Y}^- = \{y_{C+1}, \dots, y_{C+M}\}$ for classification:

$$S(x) = \frac{\text{sim}(x, \mathcal{Y})}{\text{sim}(x, \mathcal{Y}) + \text{sim}(x, \mathcal{Y}^-)} \quad (2)$$

where $\text{sim}(x, \cdot)$ represents the sum of CLIP similarities between the sample and the labels in a given label set.

3.2. Overview of SynOOD

Our proposed SynOOD, illustrated in Fig. 2, addresses this issue through a three-step process: **1) Near-Boundary OOD Image Generation.** In Fig. 2(a), an MLLM is employed to generate multiple semantic labels for each element within an image, excluding the main object. A novel iterative generative approach utilizing a diffusion model is then applied to generate near-boundary OOD images. These synthetic images help us to fine-tune the CLIP models effectively. **2) Fine-tuning of the CLIP image encoder.** We train a projection layer following the CLIP image encoder using both InD data and synthetic OOD images along with the negative labels. The image encoder remains frozen, while only the projection layer updates. **3) Fine-tuning of the CLIP text encoder features.** We make the features (the output of the CLIP text encoder) associated with a subset of

negative labels related to synthetic OOD images learnable and fine-tune them with synthetic images. This step reduces the semantic gap between InD and negative labels to an appropriate distance, improving image-text alignment.

We fine-tune the CLIP image encoder and text encoder features separately to maintain training stability. Experiments in Table 5 validate the effectiveness of this approach. After the fine-tuning of the CLIP encoders, our approach boosts the performance of OOD detection.

3.3. Near-boundary OOD Image Generation

In this image-generation process, we need to use three models: an MLLM, an in-painting diffusion model, and a traditional recognition model as the OOD detection model. Initially, we employ the MLLM ϕ to generate multiple semantic labels for each element in every InD image x^{in} , excluding the main object. The output, denoted as p^{con} is obtained as follows:

$$p^{\text{con}} = \phi(x^{\text{in}}, p^{\text{in}}), \quad (3)$$

where p^{in} is the input prompt for ϕ . Rather than relying on masks, we implement an iterative generative process when employing an in-painting diffusion model set to a strength of less than 1, enabling the generation of OOD images with minimal manual intervention.

Concretely, x^{in} and p^{con} are fed in the in-painting diffusion model to generate an image x^{syn} . Specifically, we denote the feature of x^{in} as z^{in} after passing it through the VAE [23] encoder f . Given a learnable random noise ϵ , a variance schedule $\{\alpha_1, \dots, \alpha_T\}$, a timestep T , and z^{in} , the diffusing process χ can be expressed as:

$$\begin{aligned} z_T &= \chi(z^{\text{in}}, T, \epsilon) \\ &= \sqrt{\bar{\alpha}_T} z^{\text{in}} + \sqrt{1 - \bar{\alpha}_T} \epsilon, \epsilon \sim \mathcal{N}(0, I). \end{aligned} \quad (4)$$

where $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$. In the denoising process, a U-Net [44], denote as ϵ_θ , is utilized to predict a noise needed to reconstruct z_{t-1} from z_t using a text prompt p^{con} and a mask M :

$$\begin{aligned} z_{t-1} &= \sqrt{\alpha_{t-1}} \left(\frac{z_t - \sqrt{1 - \alpha_t} \epsilon_\theta(z_t, t, P, M)}{\sqrt{\alpha_t}} \right) \\ &\quad + \sqrt{1 - \alpha_{t-1}} \epsilon_\theta(z_t, t, P, M), \end{aligned} \quad (5)$$

where $P = \psi(p^{\text{con}})$ stands for the feature of p^{con} extracted by a text encoder ψ , and M is initialized as a matrix filled with ones. By iterating through multiple time steps, the image features are gradually denoised and completed until the image at $t = 0$ is completely generated. The synthetic image is obtained through the complete denoising process, represented by γ :

$$x^{\text{syn}} = h(\gamma(z_T, T, P, M)). \quad (6)$$

where h is the VAE decoder.

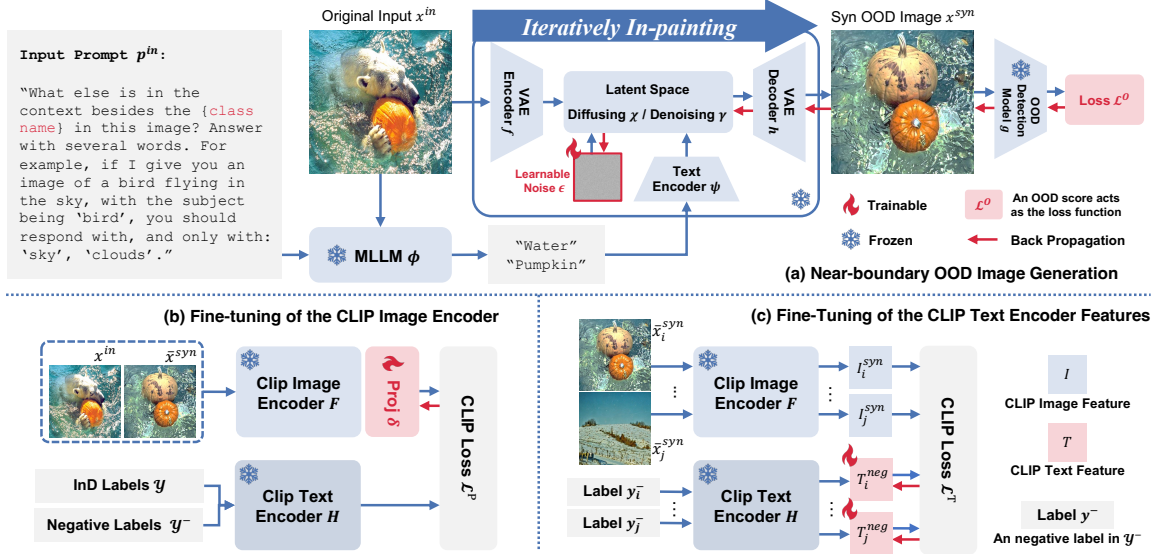


Figure 2. Overview of our proposed SynOOD framework. (a) Near-boundary OOD image generation: utilizes an MLLM and a diffusion model to iteratively generate synthetic OOD images from InD images, guided by an OOD score as the loss function. (b) Fine-tuning of the CLIP image encoder: trains a projection to strengthen connections between challenging OOD samples and negative labels. (c) Fine-tuning of the CLIP text encoder features: refines the negative label features derived from CLIP to improve OOD discrimination further.

We employ an off-the-shelf OOD detection method such as Energy score [31] as a loss function on a traditional recognition model g (e.g. ResNet50 [14]). The loss function is defined as:

$$\mathcal{L}^O = m_{\text{out}} - \tau \cdot \log \sum_{i=1}^C e^{g_i(x^{\text{syn}})/\tau}, \quad (7)$$

where m_{out} is a constant representing the OOD threshold of the OOD score, as used in [31], τ is a temperature parameter, and $g_i(x)$ denotes the logits of g for the i -th class among C categories. Combining the Eqs. (4), (6), and (7), we can calculate the gradient of the random noise ϵ at the very beginning:

$$\nabla_{\epsilon} \mathcal{L}^O = \frac{\partial \mathcal{L}}{\partial x^{\text{syn}}} \cdot \frac{\partial g}{\partial \gamma} \cdot \frac{\partial \gamma}{\partial z_T} \cdot \frac{\partial z_T}{\partial \epsilon}. \quad (8)$$

By iterative refining ϵ for a few iterations, we observe rapid convergence of the loss function, leading to the generation of highly reliable near-boundary OOD images \bar{x}^{syn} .

While calculating the gradient of ϵ can be computationally demanding, we address this challenge by adopting the Skip Gradient operation proposed by Chen et al. [5]:

$$\nabla_{\epsilon} \mathcal{L}^O \approx \ddot{\nabla}_{\epsilon} \mathcal{L}^O = \rho \cdot \frac{\partial \mathcal{L}}{\partial x^{\text{syn}}} \cdot \frac{\partial g}{\partial \gamma} \quad (9)$$

This technique significantly reduces the computational burden, enabling more efficient training. The noise updating equation can be expressed as:

$$\epsilon := \epsilon - r \cdot \ddot{\nabla}_{\epsilon} \mathcal{L}^O, \quad (10)$$

where r is the learning rate.

3.4. Fine-tuning of the CLIP image encoder

In this section, we fine-tune the CLIP image encoder by utilizing our dataset \mathcal{D}^{pro} . Specifically, after completing the generation process, we acquire a set of synthetic OOD images, denoted as \mathcal{X}^{syn} , which is paired with corresponding InD data. Each image in \mathcal{X}^{syn} is fed into the CLIP model along with the negative label set \mathcal{Y}^- , aligning a negative label to each synthetic image and forming the dataset $\mathcal{D}^{\text{syn}} = \{(x^{\text{syn}}, y^-)\}$, $x^{\text{syn}} \in \mathcal{X}^{\text{syn}}$, $y^- \in \mathcal{Y}_*^-$, where \mathcal{Y}_*^- is the subset of negative labels associated with these images. Each generated OOD image is paired one-to-one with a corresponding InD image.

Using both InD data and these synthetic OOD samples, we create a training dataset $\mathcal{D}^{\text{pro}} = \mathcal{D}^{\text{syn}} \cup \mathcal{D}_*^{\text{in}} = \{(x_i, y_i)\}_{i=1}^{2m}$ for the image encoder fine-tuning, where $\mathcal{D}_*^{\text{in}}$ represents a subset of \mathcal{D}^{in} , and m stands for the number of InD data we selected from \mathcal{D}^{in} . The selection of InD data is critical, as it directly affects the fine-tuning outcomes of the CLIP image encoder. To identify the most information-rich images within each category, we calculated the ratio of JPEG file size to the number of pixels for all images, sorted them accordingly, and then selected a specified number of top-ranked images from each category. This strategy ensures that we choose a batch of images with the highest complexity in each category.

As Fig. 2(b) shows, The parameters of the CLIP image encoder F remain frozen during training, and only the parameters of the projection layer δ are updated. We employ

the CLIP loss \mathcal{L}^P to train δ :

$$\hat{I}_i = \delta(F(x_i)), T_i = H(y_i^-), \quad (x_i, y_i) \in \mathcal{D}^{pro}, \quad (11)$$

$$\mathcal{L}^P = -\frac{1}{2m} \sum_{i=1}^{2m} \log \frac{\exp(\text{sim}(\hat{I}_i, T_i)/\tau)}{\sum_{j=1}^{M'} \exp(\text{sim}(\hat{I}_i, T_j)/\tau)}, \quad (12)$$

where $0 < M' \leq M$ represents the number of negative labels in the subset, H stands for CLIP text encoder, \hat{I}_i and T_i stand for image feature and text feature, respectively, and τ is the temperature parameter.

3.5. Fine-tuning of the CLIP text encoder features

The primary motivation for fine-tuning negative label features derived from the CLIP text encoder H is to ensure that the semantic representations of negative labels adapt specifically to OOD data with subtle variations from InD. While the fine-tuned CLIP image encoder aligns images to the corresponding negative labels, it may not fully capture the nuances of the synthetic OOD samples without some adaptation on the text side. This fine-tuning dynamically adjusts these representations for better generalization and reduces overfitting risks of the image encoder fine-tuning on a limited set of synthetic OOD data.

Specifically, we utilize CLIP and make a subset of the negative label, \mathcal{Y}_*^- , associated with synthetic OOD samples \mathcal{D}^{syn} , learnable. We denote the CLIP text features of \mathcal{Y}_*^- as $\mathcal{T}_*^{\text{neg}} = \{T_i^{\text{neg}}\}_{i=1}^{M'}$, where $M' \approx \frac{1}{2}M$, leaving the remaining labels in \mathcal{Y}^- fixed. This subset negative label fine-tuning reduces the semantic gap between InD and negative labels while maintaining model robustness, enabling precise detection of near-boundary OOD samples.

During fine-tuning, the learnable features $\mathcal{T}_*^{\text{neg}}$ are adjusted based on synthetic OOD images \mathcal{X}^{syn} , allowing the negative label features to move closer in feature space to these OOD samples while maintaining separation from InD representations. The objective is to align negative labels to capture relevant distinctions from InD data without overlap. This direct fine-tuning approach reduces the computational cost of modifying text encoder embeddings. The loss function \mathcal{L}^T for fine-tuning negative features derived from text encoder H is:

$$I_i^{\text{syn}} = F(x_i^{\text{syn}}), \quad (13)$$

$$\mathcal{L}^T = -\frac{1}{m} \sum_{i=1}^m \log \frac{\exp(\text{sim}(I_i, T_i^{\text{neg}})/\tau)}{\sum_{j=1}^{M'} \exp(\text{sim}(T_i^{\text{neg}}, I_j)/\tau)}. \quad (14)$$

Through fine-tuning, the negative text encoder features are adjusted better to capture the distinctions between InD and OOD data. We do not use the fine-tuned CLIP image encoder, as we aim to avoid adapting the image feature projection based on the text features. This approach helps prevent the negative label features from shifting to a suboptimal position, ensuring better control over their alignment. We will

provide a more detailed discussion in the experiments section.

4. Experiments

4.1. Experimental Setup

Dataset. We follow Huang et al. [20] and conduct extensive experiments with the standard large-scale ImageNet-1k [45] as InD data. For OOD data, we employ iNaturalist [51], SUN [56], Places365 [64], and Texture [6]. Moreover, we test our SynOOD on OpenOOD benchmark [58, 60]. Specifically, ImageNet-O [17], SSB-hard [52] and NINCO [2] are labeled as near-OOD, and far-OOD contains iNaturalist [51], Texture [6], and OpenImage-O [53].

Computational Cost. Compared to NegLabel, our method adds less than 1% additional parameters and takes under 2 ms per image during inference.

Implementation Details. We use LLaVA [30] to generate prompts for the diffusion model and employ Stable Diffusion 2 for inpainting [43] as the generative model to create near-boundary OOD images. We set the strength parameter to 0.6 and the number of timesteps to 20. Energy [31] is utilized as the OOD Loss function, with ResNet50 serving as the backbone model. The inpainting process iterates 3 times, with $r \cdot \rho = 10$. For the CLIP image encoder fine-tuning, we only train for 3 epochs using Adam with a learning rate of 1×10^{-3} , a batch size of 128, and a weight decay of 1×10^{-5} . For the CLIP text encoder features fine-tuning, we employ SGD with a learning rate of 2×10^{-3} and train for 5 epochs. All experiments are performed using PyTorch [39] on two NVIDIA V100.

4.2. Main Results

As presented in Table 1, we evaluate SynOOD against a range of existing OOD detection approaches on the widely-used ImageNet-1k benchmark, showcasing its performance across multiple challenging datasets. The methods listed from MSP [15] to ReAct [47] represent OOD detection approaches based on single-modal vision networks, while the methods from ZOC [11] to NegLabel [22] employ the multi-modal capabilities of the CLIP model. The results consistently demonstrate that pre-trained multi-modal models like CLIP have a significant advantage over traditional single-modal vision networks for OOD detection, underscoring the benefits of aligning both text and visual representations to improve OOD performance. When comparing SynOOD to NegLabel, we observe that SynOOD maintains strong detection performance on the iNaturalist [51] and SUN [56] datasets, which involve natural images and complex scenes, respectively. More notably, SynOOD achieves significant improvements on the Places [64] and Texture [6] datasets, which feature a broader diversity of environmental and textural variations. These improvements underscore

Table 1. Comparison of OOD detection performance between SynOOD and existing methods. The best and second-best results are highlighted in **bold** and underlined, respectively. All methods use ViT/B-16 as the backbone. Methods in the upper section are pre-trained on ImageNet, while those in the lower section utilize CLIP pre-training.

Method	OOD Datasets								Average	
	iNaturalist		SUN		Place		Texture			
	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓
MSP [15]	87.44	58.36	79.73	73.72	79.67	74.41	79.69	71.93	81.63	69.61
ODIN [28]	94.65	30.22	87.17	54.04	85.54	55.06	87.85	51.67	88.80	47.75
Energy [31]	95.33	26.12	92.66	35.97	91.41	39.87	86.76	57.61	91.54	39.89
GradNorm [21]	72.56	81.50	72.86	82.00	73.70	80.41	70.26	79.36	72.35	80.82
ViM [53]	93.16	32.19	87.19	54.01	83.75	60.67	87.18	53.94	87.82	50.20
KNN [48]	94.52	29.17	92.67	35.62	91.02	39.61	85.67	64.35	90.97	42.19
VOS [9]	94.62	28.99	92.57	36.88	91.23	38.39	86.33	61.02	91.19	41.32
DICE [46]	94.49	25.63	90.83	35.15	87.48	46.49	90.30	31.72	90.78	34.75
ReAct [47]	96.22	20.38	94.20	24.20	91.58	33.85	89.80	47.30	92.95	31.43
ZOC [11]	86.09	87.30	81.20	81.51	83.39	73.06	76.46	98.90	81.79	85.19
MCM [35]	94.59	32.20	92.25	38.80	90.31	46.20	86.12	58.50	90.82	43.93
CoOp [66]	94.89	29.47	93.36	31.34	90.07	40.28	87.58	54.25	91.48	38.84
CoCoOp [65]	94.73	30.74	93.15	31.18	90.63	38.75	87.92	53.84	91.61	38.63
NPOS [49]	96.19	16.58	90.44	43.77	89.44	45.27	88.80	46.12	91.22	37.94
CLIPN [54]	95.27	23.94	93.93	26.17	92.28	33.45	90.93	40.83	93.10	31.10
LSN [37]	95.83	21.56	94.35	26.32	91.25	34.48	90.42	38.54	92.96	30.23
LoCoOp [36]	96.86	16.05	95.07	23.44	91.98	32.87	90.19	42.18	93.53	28.64
NegLabel [22]	99.49	1.91	95.49	20.53	91.64	35.59	90.22	43.56	94.21	25.40
CSP[3]	<u>99.60</u>	<u>1.54</u>	<u>96.66</u>	<u>13.66</u>	92.90	<u>29.32</u>	93.86	<u>25.52</u>	95.76	<u>17.51</u>
AdaNeg[61]	99.71	0.59	97.44	9.50	94.55	34.34	94.93	31.27	96.66	18.92
SynOOD (Ours)	99.57	1.57	95.82	20.46	97.37	12.12	95.29	22.94	97.01	14.27

Table 2. OOD detection performance on the OpenOOD benchmark. The methods in the upper section are using the whole ImageNet for training. Results are averaged across OOD datasets.

Method	AUROC↑		FPR95↓	
	NearOOD	FarOOD	NearOOD	FarOOD
GEN [32]	78.97	90.98	-	-
AugMix [16]+ReAct [47]	79.94	93.70	-	-
RMDS [42]	80.09	92.60	-	-
AugMix [16]+ASH [7]	82.16	96.05	-	-
MCM [35]	59.89	80.71	81.02	68.88
NegLabel [22]	75.47	94.30	74.74	25.73
Ours	77.55	96.21	71.68	17.11

SynOOD’s ability to more accurately capture and represent OOD boundaries in data with high intra-class variability and complex visual patterns, areas where traditional methods often struggle. Overall, SynOOD establishes a new state-of-the-art in OOD detection, surpassing previous methods with a substantial AUROC improvement of 2.80% and an FPR95 reduction of 11.13%. This performance boost reflects SynOOD’s robust design and effective use of negative label fine-tuning and iterative OOD sample generation, which together enable a more nuanced alignment between InD and near-boundary OOD samples.

Evaluation on OpenOOD benchmark. We further evaluate SynOOD on the OpenOOD benchmark, which includes

both near-OOD and far-OOD scenarios. As presented in Table 2, the methods in the upper section are drawn from OpenOOD [60]. These methods typically show stronger performance in near-OOD detection, as they benefit from training on the full ImageNet dataset, which contains over 1.2 million images, giving them a substantial advantage. This allows them to capture diverse InD patterns, improving near-OOD detection accuracy. In contrast, SynOOD uses only a lightweight subset of 50k ImageNet images yet achieves competitive performance, outperforming all methods in far-OOD detection and exceeding MCM [35] and NegLabel [22] on near-OOD detection. These findings underscore SynOOD’s effectiveness across both near and far-OOD conditions, even with limited training data. This balance highlights the robustness of our approach, particularly in the challenging far-OOD scenario, where our method consistently maintains superior discrimination. These results confirm the generalization capability of SynOOD and its adaptability across various OOD conditions.

4.3. Ablation Study

Image Generation and Training Components. In Tab. 4, we investigate the effect of the fine-tuning of the CLIP image encoder and the fine-tuning of the CLIP text encoder features using various synthetic image generation strategies. We compare two image generation approaches for synthesizing OOD data: (1) directly generating images using neg-

Table 3. OOD detection performance comparison across various CLIP architectures. Results are averaged across four OOD datasets.

Backbone	Method	OOD Datasets								Average	
		iNatrualist		SUN		Place		Texture		AUROC↑	FPR95↓
		AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓		
ResNet50	NegLabel	99.24	2.88	94.54	26.51	89.72	42.60	88.40	50.8	92.97	30.70
	SynOOD	98.89	4.26	95.13	24.98	96.07	16.29	93.09	34.65	95.80	20.05
ViT-B/32	NegLabel	99.11	3.73	95.27	22.48	91.92	34.94	88.57	50.51	93.67	27.92
	SynOOD	99.35	2.42	94.90	25.10	97.20	13.24	92.26	36.86	95.93	19.41
ViT-B/16	NegLabel	99.49	1.91	95.49	20.53	91.64	35.59	90.22	43.56	94.21	25.40
	SynOOD	99.57	1.57	95.82	20.46	97.37	12.12	95.29	22.94	97.01	14.27

Table 4. A set of ablation experiments on SynOOD. “FT Label” refers to fine-tuning the label features, “Neg Image” indicates images generated by a text-to-image diffusion model prompted by the negative labels, and “Grad Image” represents images produced using our iterative generation process. Results are averaged across four OOD datasets.

Projection	FT Label	Data		Average	
		Neg Image	Grad Image	AUROC↑	FPR95↓
-	-	-	-	94.21	25.40
✓	-	✓	-	94.85	22.59
✓	-	-	✓	96.21	17.29
✓	✓	✓	-	95.01	21.93
✓	✓	✓	✓	97.01	14.27

ative labels as prompts in a text-to-image diffusion model and (2) using an iterative image generation process that refines the alignment between generated images and their associated negative labels. The first method, direct generation, employs negative labels as prompts to synthesize images in a single pass through a diffusion model. While effective, this approach can sometimes yield images with limited similarity with InD data, which may not fully capture the nuanced distinctions between InD and near-boundary OOD samples. The iterative approach substantially improves performance across both image encoder fine-tuning and label feature fine-tuning, as it progressively shifts the image away from the original theme while still preserving visual connections to the InD data through the background features. Furthermore, our experiments indicate that the best average performance is achieved when both the image encoder fine-tuning and negative label fine-tuning are jointly applied with the iterative generation strategy.

Effect of the number of synthetic data. Figure 3 shows how varying the amount of synthetic OOD data affects SynOOD when fine-tuning both the CLIP image encoder and negative text encoder features. We evaluate SynOOD with $m \in \{1k, 5k, 10k, 20k, 30k, 50k, 75k, 100k\}$ synthetic OOD samples. SynOOD maintains stable performance, with accuracy improving as m increases, demonstrating the effectiveness of our iterative data generation

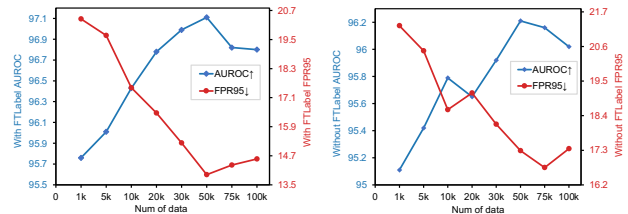


Figure 3. SynOOD performance with different amounts of synthetic OOD data. The left plot shows results with fine-tuning negative label features, while the right plot shows results without fine-tuning. Results are averaged across four OOD datasets.

strategy. However, performance slightly declines when m exceeds 50k, due to the balance between fixed and fine-tuned negative label features. Our method fine-tunes about 5.5k of 11k negative labels, preserving strong performance on far-boundary OOD data and enhancing sensitivity to near-boundary samples. Fine-tuning over 7k labels (e.g., $m = 75k$) disrupts this balance, leading to misclassification of easier OOD samples and a minor performance drop. Additionally, including fine-tuned negative label features consistently improves performance across all m values compared to training without fine-tuning, underscoring the importance of tailored feature alignment for effective OOD detection.

Analysis of different CLIP networks. In Table 3, we evaluate the effectiveness of our method, SynOOD, across a range of CLIP-based architectures, including ResNet50 [14], ViT-B/32 [8], and ViT-B/16 [8]. SynOOD consistently outperforms the baseline, NegLabel [22], demonstrating superior performance across all architectures. Notably, our method yields significant improvements in the FPR95 metric, achieving reductions exceeding 10% on each network, which is a substantial enhancement for high-confidence OOD detection. This performance boost indicates that SynOOD is not only effective in lowering false positive rates but also exhibits strong robustness and adaptability across diverse backbone architectures. The consistent gains achieved across both convolutional (ResNet) and transformer-based (ViT) models un-

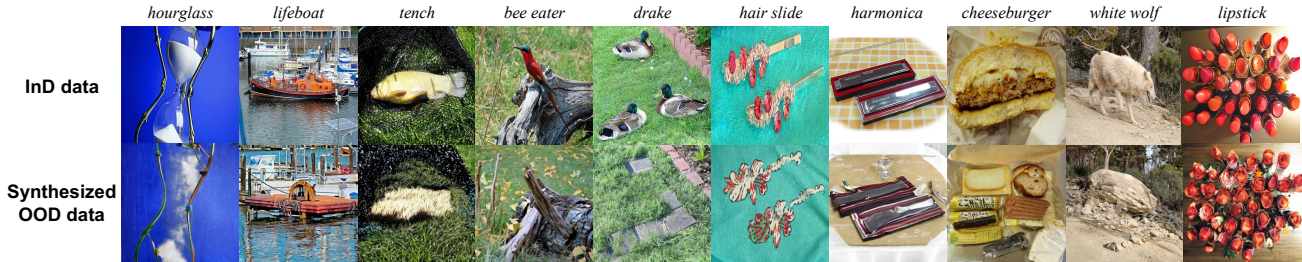


Figure 4. Visualization of InD and synthetic OOD data. The labels at the top of the figure represent ImageNet categories as InD. The first row shows ImageNet (InD) images, while the second row presents our synthetic OOD data.

Table 5. Performance comparison across training strategies, including joint training and step-by-step training with or without the trained projection layer during fine-tuning. Results are averaged across four OOD datasets.

Strategy	Fine-Tuning with Projection	Average	
		AUROC↑	FPR95↓
Joint	-	95.93	19.41
Separate	✓	96.49	16.18
Separate	✗	97.01	14.27

derscore the generalizability of our approach, showing that SynOOD’s design principles are broadly applicable to various network structures within the CLIP framework. This adaptability further emphasizes SynOOD’s potential for application in a wide range of OOD detection scenarios.

Training Strategy Comparison. In Table 5, we examine three distinct training strategies to evaluate our SynOOD framework: joint training and two step-by-step training approaches. For joint training, we optimize both the projection layer and label feature fine-tuning. The step-by-step approach separates the training into sequential phases. Specifically, we implement two variations of step-by-step training: one that incorporates the pre-trained projection layer during the label feature fine-tuning phase, and another that excludes it. The results in Table 5 reveal that step-by-step training is more effective than joint training, providing both enhanced stability during training and improved detection performance. This improvement is due to the sequential focus on each component, which may reduce interference effects seen in joint training. Notably, our experiments indicate that fine-tuning the label features without image encoder fine-tuning yields the best results. We hypothesize that the projection layer, trained on only 50k synthetic OOD samples, is prone to overfitting, which may limit generalization when combined with fine-tuned label features. By excluding the projection layer in this phase, SynOOD achieves a better balance between specificity and robustness in OOD detection. This evaluation of training strategies underscores the importance of carefully structuring the training process

for complex OOD detection systems.

Visualize of the synthetic Data. In Fig. 4, we present several InD and synthetic OOD data pairs to illustrate how our method generates OOD samples that are visually similar to InD samples, yet exhibit clear OOD characteristics. These examples show that the synthetic OOD images closely resemble their corresponding InD images to the human eye, but contain subtle differences that distinguish them as OOD. For instance, in the hourglass image on the left of Fig. 4, the original InD sample depicts an hourglass with white sand, slim supports, and a blue background. Our generation process has transformed this scene into an image with a blue sky, white clouds, and vines, making it perceptually similar yet meaningfully different from the InD data. Similarly, on the right, lipsticks are reimagined as flowers in the same setting. Our method can even decompose a cheeseburger, displaying each component separately against a consistent background. We hope our iterative generation process inspires further research into innovative OOD data generation techniques and opens new possibilities for other applications where similar approaches might be beneficial.

5. Conclusion

In this paper, we present SynOOD, a novel approach to OOD detection that combines iterative generative techniques with fine-tuned CLIP models and features for enhanced identification of challenging OOD samples. By synthesizing near-boundary OOD samples using a diffusion model, SynOOD generates data that is visually similar to, yet semantically distinct from, InD data, allowing for more precise OOD discrimination. Extensive evaluations on multiple benchmark datasets demonstrate that SynOOD surpasses existing methods, achieving state-of-the-art performance in AUROC and FPR95 metrics. Our work highlights the effectiveness of synthetic data for OOD detection, suggesting new directions for using generative methods to improve model robustness in diverse tasks. We believe SynOOD opens up new directions for OOD detection and encourages future research into similar generative strategies for improving model robustness across other tasks.

6. Acknowledgement

This work was supported by National Natural Science Foundation of China (No.62072112) and Shanghai Science and Technology Committee under Grant (No. 24511103900, 24511103202) and was partly supported by National Key RD Program of China under grant No. 2022YFC3601405.

References

- [1] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023. 1
- [2] Julian Bitterwolf, Maximilian Mueller, and Matthias Hein. In or out? fixing imagenet out-of-distribution detection evaluation. *arXiv preprint arXiv:2306.00826*, 2023. 5
- [3] Mengyuan Chen, Junyu Gao, and Changsheng Xu. Conjugated semantic pool improves ood detection with pre-trained vision-language models. *Advances in Neural Information Processing Systems*, 37:82560–82593, 2024. 6
- [4] Xingyu Chen, Xuguang Lan, Fuchun Sun, and Nanning Zheng. A boundary based out-of-distribution classifier for generalized zero-shot learning. In *European conference on computer vision*, pages 572–588. Springer, 2020. 2
- [5] Zhaoyu Chen, Bo Li, Shuang Wu, Kaixun Jiang, Shouhong Ding, and Wenqiang Zhang. Content-based unrestricted adversarial attack. In *Advances in Neural Information Processing Systems*, pages 51719–51733. Curran Associates, Inc., 2023. 4
- [6] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3606–3613, 2014. 5
- [7] Andrija Djurisic, Nebojsa Bozanic, Arjun Ashok, and Rosanne Liu. Extremely simple activation shaping for out-of-distribution detection. *arXiv preprint arXiv:2209.09858*, 2022. 6
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021. 1, 7
- [9] Xuefeng Du, Zhaoning Wang, Mu Cai, and Yixuan Li. Vos: Learning what you don't know by virtual outlier synthesis. *arXiv preprint arXiv:2202.01197*, 2022. 6
- [10] Xuefeng Du, Yiyu Sun, Jerry Zhu, and Yixuan Li. Dream the impossible: Outlier imagination with diffusion models. *Advances in Neural Information Processing Systems*, 36:60878–60901, 2023. 3
- [11] Sepideh Esmaeilpour, Bing Liu, Eric Robertson, and Lei Shu. Zero-shot out-of-distribution detection based on the pre-trained model clip. In *Proceedings of the AAAI conference on artificial intelligence*, pages 6568–6576, 2022. 5, 6
- [12] Zhen Fang, Yixuan Li, Jie Lu, Jiahua Dong, Bo Han, and Feng Liu. Is out-of-distribution detection learnable? *Advances in Neural Information Processing Systems*, 35:37199–37213, 2022. 2
- [13] Stanislav Fort, Jie Ren, and Balaji Lakshminarayanan. Exploring the limits of out-of-distribution detection. *Advances in Neural Information Processing Systems*, 34:7068–7081, 2021. 2
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4, 7
- [15] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016. 1, 2, 5, 6
- [16] Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. Augmix: A simple data processing method to improve robustness and uncertainty. *arXiv preprint arXiv:1912.02781*, 2019. 6
- [17] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15262–15271, 2021. 5
- [18] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017. 1
- [19] Haiwen Huang, Zhihan Li, Lulu Wang, Sishuo Chen, Bin Dong, and Xinyu Zhou. Feature space singularity for out-of-distribution detection. *arXiv preprint arXiv:2011.14654*, 2020. 2
- [20] Rui Huang and Yixuan Li. Mos: Towards scaling out-of-distribution detection for large semantic space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8710–8719, 2021. 5
- [21] Rui Huang, Andrew Geng, and Yixuan Li. On the importance of gradients for detecting distributional shifts in the wild. *Advances in Neural Information Processing Systems*, 34:677–689, 2021. 1, 2, 6
- [22] Xue Jiang, Feng Liu, Zhen Fang, Hong Chen, Tongliang Liu, Feng Zheng, and Bo Han. Negative label guided ood detection with pretrained vision-language models. *arXiv preprint arXiv:2403.20078*, 2024. 1, 3, 5, 6, 7
- [23] Diederik P Kingma. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 3
- [24] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *Advances in neural information processing systems*, 31, 2018. 2

- [25] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 1
- [26] Jinglun Li, Xinyu Zhou, Pinxue Guo, Yixuan Sun, Yiwen Huang, Weifeng Ge, and Wenqiang Zhang. Hierarchical visual categories modeling: A joint representation learning and density estimation framework for out-of-distribution detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 23425–23435, 2023. 1
- [27] Jinglun Li, Xinyu Zhou, Kaixun Jiang, Lingyi Hong, Pinxue Guo, Zhaoyu Chen, Weifeng Ge, and Wenqiang Zhang. Tagood: A novel approach to out-of-distribution detection via vision-language representations and class center learning. *arXiv preprint arXiv:2408.15566*, 2024. 1
- [28] Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. *arXiv preprint arXiv:1706.02690*, 2017. 1, 2, 6
- [29] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. 1
- [30] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 5
- [31] Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. *Advances in neural information processing systems*, 33:21464–21475, 2020. 1, 4, 5, 6
- [32] Xixi Liu, Yaroslava Lochman, and Christopher Zach. Gen: Pushing the limits of softmax-based out-of-distribution detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23946–23955, 2023. 6
- [33] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. 1
- [34] George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995. 3
- [35] Yifei Ming, Ziyang Cai, Jiuxiang Gu, Yiyou Sun, Wei Li, and Yixuan Li. Delving into out-of-distribution detection with vision-language representations. *Advances in neural information processing systems*, 35:35087–35102, 2022. 1, 2, 6
- [36] Atsuyuki Miyai, Qing Yu, Go Irie, and Kiyoharu Aizawa. Locoop: Few-shot out-of-distribution detection via prompt learning. *Advances in Neural Information Processing Systems*, 36, 2024. 1, 6
- [37] Jun Nie, Yonggang Zhang, Zhen Fang, Tongliang Liu, Bo Han, and Xinmei Tian. Out-of-distribution detection with negative prompts. In *The Twelfth International Conference on Learning Representations*, 2024. 3, 6
- [38] OpenAI. Gpt-4 technical report, 2023. Accessed: 2023-10-23. 1
- [39] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. 5
- [40] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution image synthesis. In *The Twelfth International Conference on Learning Representations*, 2024. 1
- [41] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 1
- [42] Jie Ren, Stanislav Fort, Jeremiah Liu, Abhijit Guha Roy, Shreyas Padhy, and Balaji Lakshminarayanan. A simple fix to mahalanobis distance for improving near-ood detection. *arXiv preprint arXiv:2106.09022*, 2021. 6
- [43] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1, 5
- [44] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015. 3
- [45] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252, 2015. 5
- [46] Yiyou Sun and Yixuan Li. Dice: Leveraging sparsification for out-of-distribution detection. In *European Conference on Computer Vision*, pages 691–708. Springer, 2022. 1, 6
- [47] Yiyou Sun, Chuan Guo, and Yixuan Li. React: Out-of-distribution detection with rectified activations. *Advances in Neural Information Processing Systems*, 34:144–157, 2021. 1, 2, 5, 6
- [48] Yiyou Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. Out-of-distribution detection with deep nearest neighbors. In *International Conference on Machine Learning*, pages 20827–20840. PMLR, 2022. 6
- [49] Leitian Tao, Xuefeng Du, Xiaojin Zhu, and Yixuan Li. Non-parametric outlier synthesis. *arXiv preprint arXiv:2303.02966*, 2023. 3, 6
- [50] Joost Van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. Uncertainty estimation using a single deep deterministic neural network. In *International conference on machine learning*, pages 9690–9700. PMLR, 2020. 2
- [51] Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and

- Serge Belongie. The inaturalist species classification and detection dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8769–8778, 2018. 5
- [52] Sagar Vaze, Kai Han, Andrea Vedaldi, and Andrew Zisserman. Open-set recognition: A good closed-set classifier is all you need. In *International Conference on Learning Representations*, 2021. 5
- [53] Haoqi Wang, Zhizhong Li, Litong Feng, and Wayne Zhang. Vim: Out-of-distribution with virtual-logit matching. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4921–4930, 2022. 5, 6
- [54] Hualiang Wang, Yi Li, Hui Feng Yao, and Xiaomeng Li. Clipn for zero-shot ood detection: Teaching clip to say no. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1802–1812, 2023. 1, 2, 6
- [55] Jim Winkens, Rudy Bunel, Abhijit Guha Roy, Robert Stanforth, Vivek Natarajan, Joseph R Ledsam, Patricia MacWilliams, Pushmeet Kohli, Alan Karthikesalingam, Simon Kohl, et al. Contrastive training for improved out-of-distribution detection. *arXiv preprint arXiv:2007.05566*, 2020. 2
- [56] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *2010 IEEE computer society conference on computer vision and pattern recognition*, pages 3485–3492. IEEE, 2010. 5
- [57] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017. 1
- [58] Jing Kang Yang, Pengyun Wang, Dejian Zou, Zitang Zhou, Kunyuan Ding, Wenxuan Peng, Haoqi Wang, Guangyao Chen, Bo Li, Yiyu Sun, et al. Openood: Benchmarking generalized out-of-distribution detection. *Advances in Neural Information Processing Systems*, 35:32598–32611, 2022. 5
- [59] Alireza Zaemzadeh, Niccolo Bisagno, Zeno Sambugaro, Nicola Conci, Nazanin Rahnavard, and Mubarak Shah. Out-of-distribution detection using union of 1-dimensional subspaces. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, pages 9452–9461, 2021. 2
- [60] Jingyang Zhang, Jing Kang Yang, Pengyun Wang, Haoqi Wang, Yueqian Lin, Haoran Zhang, Yiyu Sun, Xuefeng Du, Kaiyang Zhou, Wayne Zhang, et al. Openood v1. 5: Enhanced benchmark for out-of-distribution detection. *arXiv preprint arXiv:2306.09301*, 2023. 5, 6
- [61] Yabin Zhang and Lei Zhang. Adaneg: Adaptive negative proxy guided OOD detection with vision-language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 6
- [62] Yabin Zhang, Wenjie Zhu, Chenhang He, and Lei Zhang. Lapt: Label-driven automated prompt tuning for ood detection with vision-language models. In *European Conference on Computer Vision*, pages 271–288. Springer, 2024. 3
- [63] Haotian Zheng, Qizhou Wang, Zhen Fang, Xiaobo Xia, Feng Liu, Tongliang Liu, and Bo Han. Out-of-distribution detection learning with unreliable out-of-distribution sources. *Advances in neural information processing systems*, 36:72110–72123, 2023. 2
- [64] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE transactions on pattern analysis and machine intelligence*, 40(6):1452–1464, 2017. 5
- [65] Kaiyang Zhou, Jing Kang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16816–16825, 2022. 6
- [66] Kaiyang Zhou, Jing Kang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022. 6