

# Adaptive Prompt Learning via Gaussian Outlier Synthesis for Out-of-distribution Detection

Anonymous ICCV submission

Paper ID 15742

## Abstract

001 *Out-of-distribution (OOD) detection aims to distinguish*  
 002 *whether detected objects belong to known categories or not.*  
 003 *Existing methods extract OOD samples from In-distribution*  
 004 *(ID) data to regularize the model’s decision boundaries.*  
 005 *However, the decision boundaries are not adequately regu-*  
 006 *larized due to the model’s lack of knowledge about the dis-*  
 007 *tribution of OOD data. To address the above issue, we pro-*  
 008 *pose an Adaptive Prompt Learning framework via Gaussian*  
 009 *Outlier Synthesis (APLGOS) for OOD detection. Specifi-*  
 010 *cally, we leverage the Vision-Language Model (VLM) to ini-*  
 011 *tialize learnable ID prompts by sampling standardized re-*  
 012 *sults from pre-defined Q&A pairs. Region-level prompts are*  
 013 *synthesised in low-likelihood regions of class-conditional*  
 014 *gaussian distributions. These prompts are then utilized to*  
 015 *initialize learnable OOD prompts and optimized with adap-*  
 016 *tive prompt learning. Also, OOD pseudo-samples are syn-*  
 017 *thesised via gaussian outlier synthesis. Similarity score*  
 018 *between prompts and images is utilized to calculate con-*  
 019 *trastive learning loss in high-dimensional hidden space.*  
 020 *The aforementioned methodology regularizes the model to*  
 021 *learn more compact decision boundaries for ID and OOD*  
 022 *categories. Extensive experiments show that our proposed*  
 023 *method achieves state-of-the-art performance with less ID*  
 024 *data on four mainstream datasets.*

## 025 1. Introduction

026 Deep learning has made significant progress in recent years.  
 027 It encompasses a multitude of research domains, including  
 028 object detection [39, 51], autonomous driving [40, 50] and  
 029 image generation [20, 41]. Various existing deep learning  
 030 methods rely on large-scale datasets to regularize the model,  
 031 enabling it to learn sufficient data distribution and supervi-  
 032 sion signals of the training data. In real-world scenarios,  
 033 where the number of unknown categories is significantly  
 034 greater than that in the training dataset, the model lacks  
 035 knowledge about the distribution of unknown data in prac-

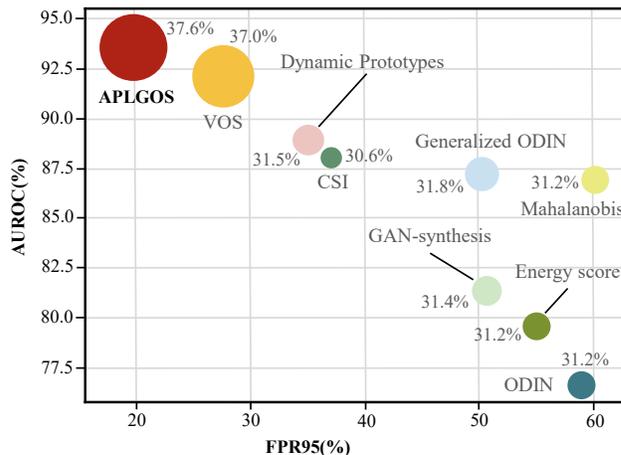


Figure 1. Quantitative comparisons with state-of-the-art OOD detection methods in terms of FPR95, AUROC and mAP metrics. Note that larger points denote higher mAP, and the numerical values are also given next to each point. Our APLGOS provides remarkable performance boost on all the metrics.

tical applications and struggles to learn compact decision boundaries that effectively distinguish between known and unknown categories. During the testing phase, unknown categories is likely to result in erroneous predictions accompanied by a high confidence score. This leads to severe safety risks in critical safety domains such as autonomous driving.

OOD detection [5, 21, 23, 32] is a research hotspot in recent years, which aims to enable the detectors to accurately distinguish not only seen categories, but also unseen categories during training. The detectors need to learn compact decision boundaries during training, ensuring low uncertainty for ID categories while maintaining high uncertainty away from them. To achieve this, existing OOD detection methods [12, 13, 33–35] provide sufficient supervision of OOD data for model training by extracting OOD pseudo-samples from ID data, helping the model better distinguish between known and unknown categories. However, due to the unpredictable quality of OOD pseudo-samples extracted from the ID data and the requirement

036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055



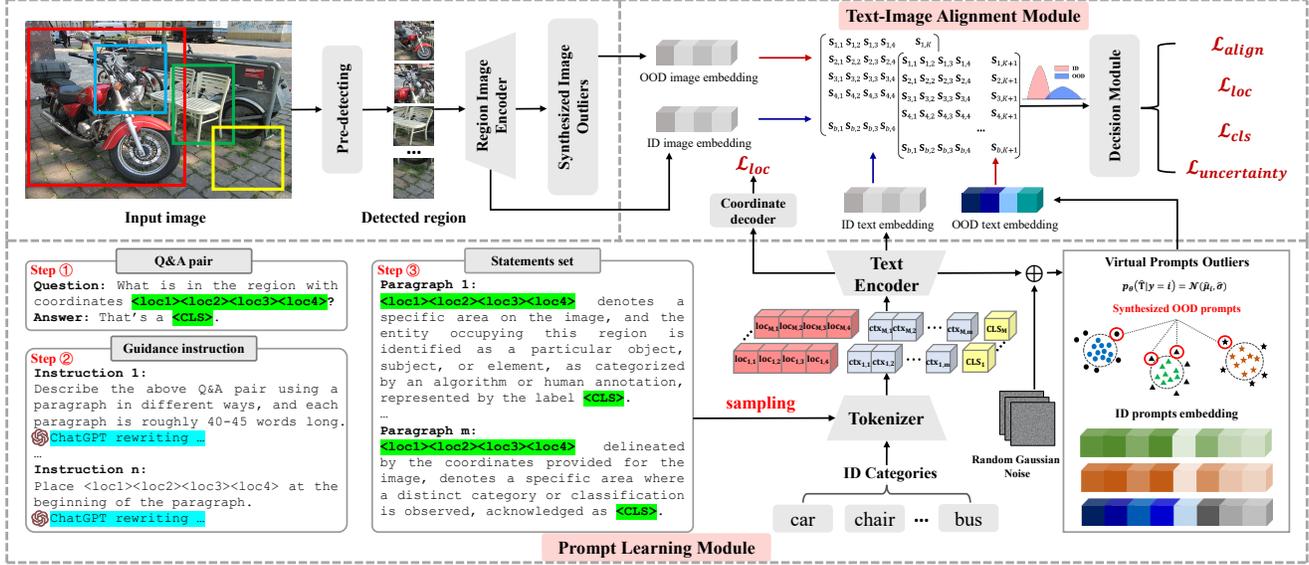


Figure 2. The proposed APLGOS network architecture. Prompt learning module is responsible for using ChatGPT to standardize Q&A pairs with guidance introduction and templates, then it generates a statements set. The module samples prompts from the statements set to initialize the learnable ID prompts, and synthesises virtual OOD prompts in low-likelihood regions of class-conditional gaussian distributions. The Text-Image Alignment Module computes similarity scores to align text and image embeddings in the hidden space.

### 157 3. Methodology

158 We propose an Adaptive Prompt Learning framework via  
 159 Gaussian Outlier Synthesis for OOD Detection. As shown  
 160 in Figure 2, APLGOS mainly consists of two modules, *i.e.*  
 161 PLM and TAM. PLM leverages ChatGPT to standardize  
 162 pre-defined Q&A pairs using guidance instructions and pre-  
 163 defined templates, generating a set of statements. During  
 164 training, PLM samples statements from this set to initialize  
 165 the learnable prompts. For ID categories, APLGOS directly  
 166 employs the initialized prompts as input to the text encoder,  
 167 whereas for OOD categories, it synthesizes virtual OOD  
 168 prompts and images within the low-likelihood region of the  
 169 class-conditional Gaussian distribution of ID classes in the  
 170 hidden space. Notably, only ID images are sourced from  
 171 the dataset, while ID prompts, OOD prompts, and OOD im-  
 172 ages are all virtual and synthesized. This approach enables  
 173 the model to enhance the quality of pseudo-samples with  
 174 less ID data while better capturing the distribution of OOD  
 175 data. Additionally, through contrastive learning, TAM com-  
 176 puts similarity scores to align images and prompts within  
 177 the high-dimensional hidden space.

178 For clarity, we omit the *batchsize* of data in the follow-  
 179 ing description and consider a single batch as an example.  
 180 The input to APLGOS consists of two modalities: detected  
 181 region images  $[X_1, X_2, \dots, X_b]$  extracted from a raw RGB  
 182 image  $X \in \mathbb{R}^{C \times H \times W}$ , and text prompts  $T \in \mathbb{R}^{b \times l}$ . Here,  
 183  $C$ ,  $H$ , and  $W$  denote the number of channels, height, and  
 184 width of the image, respectively.  $b$  represents the number  
 185 of detected region images from a single raw RGB image.

186  $l$  indicates the length of the text prompts. The text input  
 187 is given as  $T = [T_1, T_2, \dots, T_b]$ , where the  $\langle \text{CLS} \rangle$  token  
 188 in the sampled prompts has been replaced with the corre-  
 189 sponding labels.

#### 190 3.1. Prompt Learning Module

191 **ID Prompts.** To enhance the model’s representation abil-  
 192 ity and more effectively regularize its decision bound-  
 193 aries, we generate a set of statements for the Prompt  
 194 Learning Module to sample from, rather than using a  
 195 single invariant statement to initialize the learnable ID  
 196 prompts. Specifically, we first predefine a Q&A pair,  
 197 such as “*Q: What is in the region with coordinates*  
 198  *$\langle loc1 \rangle, \langle loc2 \rangle, \langle loc3 \rangle, \langle loc4 \rangle$ ? A: That’s a  $\langle \text{CLS} \rangle$ .*”  
 199 We then input this Q&A pair into ChatGPT for standardiza-  
 200 tion. During this process, we provide predefined templates  
 201 and guiding instructions to ensure that ChatGPT standard-  
 202 izes the Q&A pair accordingly. The standardization process  
 203 is illustrated below with an example prompt:

$$204 \Omega_0 = g(Q^A + M + G_0), \quad \Omega_i = g(\Omega_{i-1} + G_i), \quad (1)$$

205 where  $\Omega_i$  denotes generated prompt result in  $i_{th}$  round,  $Q^A$   
 206 denotes Q&A pair,  $M$  denotes predefined template,  $G_i$  de-  
 207 notes guidance instruction for  $i_{th}$  standardizing round and  
 208  $g$  is ChatGPT’s standardizing operation. We collect the  
 209 statements from these  $t$  rounds to obtain statements set  $\Omega_t$ .  
 210 These statements are then used for sampling during the ini-  
 211 tialization of learnable ID prompts.

212 We introduce no extra character sets and vocabularies,  
213 and the generated prompts are represented in natural lan-  
214 guage. The learnable prompts follow the paradigm *e.g.*  
215  $\langle loc_1 \rangle \langle loc_2 \rangle \langle loc_3 \rangle \langle loc_4 \rangle \langle V_1 \rangle \langle V_2 \rangle \dots \langle V_m \rangle$   
216  $\langle CLS \rangle$ , which is initialized by sampled prompt.  
217  $\langle loc_1 \rangle \langle loc_2 \rangle \langle loc_3 \rangle \langle loc_4 \rangle$  are learnable location to-  
218 kens, which implicitly introduce location information into  
219 the prompts.  $\langle V_1 \rangle \langle V_2 \rangle \dots \langle V_m \rangle$  are learnable descrip-  
220 tion tokens, and  $m$  is its length.  $\langle CLS \rangle$  is class token.

$$221 \quad \hat{\mathbf{T}} = f_\theta(h(r(g(\Omega_{t-1} + G_t))))), \quad (2)$$

222 where  $\hat{\mathbf{T}} = [\hat{\mathbf{T}}_1, \hat{\mathbf{T}}_2, \dots, \hat{\mathbf{T}}_b]$ ,  $\hat{\mathbf{T}}_i \in \mathbb{R}^{l'}$ ,  $t$  is rounds of stan-  
223 dardizing operations,  $l'$  is length of prompt embedding.  
224 Here, for ease of understanding, we use one  $\hat{\mathbf{T}}_i$  as an exam-  
225 ple to describe the subsequent operations, and standardize  
226  $\hat{\mathbf{T}}_i$  as  $\hat{\mathbf{T}}$ ,  $\hat{\mathbf{T}} \in \mathbb{R}^{l'}$ .  $f_\theta$  is transformer-based text encoder,  $h$   
227 is tokenizer,  $r$  is replacement function for  $\langle CLS \rangle$  token.  
228 We replace  $\langle CLS \rangle$  directly with the category label of the  
229 object in the current region (i.e., the corresponding ID class  
230 label).

231 **OOD Prompts.** In the hidden space, distinct decision  
232 boundaries should be established between ID and OOD  
233 prompts. In the OOD detection task, we refine the deci-  
234 sion boundaries as much as possible. By incorporating  
235 prompt learning, we synthesize region-level OOD pseudo-  
236 prompts using Gaussian outlier synthesis. Specifically, the  
237 Prompt Learning Module synthesizes virtual OOD prompts  
238 in the low-likelihood regions of class-conditional Gaussian  
239 distributions in hidden space. This allows the Text-Image  
240 Alignment Module to perceive the distribution difference  
241 between ID and OOD categories in hidden space and align  
242 images and prompts through contrastive learning. Provided  
243 that the quantity of data is large enough, we assume the  
244 ID prompts embedding from text encoder form a class-  
245 conditional multivariate Gaussian distribution:

$$246 \quad p_\theta(\hat{\mathbf{T}}|y = i) = \mathcal{N}(\hat{\mu}_i, \hat{\sigma}), \quad (3)$$

247 where  $\theta$  is the parameter of text encoder  $f_\theta$ ,  $y$  is ground truth  
248 label,  $\hat{\mu}_i$  is empirical gaussian mean of  $i_{th}$  in-distribution  
249 category prompts embedding, and  $i \in \{1, 2, \dots, K\}$ ,  $K$  rep-  
250 resents the number of in-distribution classes,  $\mathcal{N}(\hat{\mu}_i, \hat{\sigma}) =$   
251  $\frac{1}{\sqrt{2\pi\hat{\sigma}}} e^{-\frac{(\hat{\mathbf{T}} - \hat{\mu}_i)^2}{2\hat{\sigma}^2}}$ ,  $\hat{\sigma}$  denotes the tied covariance matrix.

252 First, we calculate the empirical gaussian mean of  $i_{th}$  ID  
253 category prompts embedding as follows:

$$254 \quad \hat{\mu}_i = \frac{1}{|\mathcal{Q}_T|} \sum_{j=1}^{|\mathcal{Q}_T|} \hat{\mathbf{T}}_{i,j}, \quad (4)$$

255 where  $|\mathcal{Q}_T|$  denotes the length of the prompts queue  $\mathcal{Q}_T$   
256 used to buffer ID prompts, and  $\mathcal{Q}_T \in \mathbb{R}^{K \times |\mathcal{Q}_T|}$ .

257 Then we calculate the tied covariance matrix of ID  
258 prompts embedding as follows:

$$\hat{\sigma} = \frac{1}{K|\mathcal{Q}_T|} \sum_{i=1}^K \sum_{j=1}^{|\mathcal{Q}_T|} (\hat{\mathbf{T}}_{i,j} + \alpha\varepsilon - \hat{\mu}_i)(\hat{\mathbf{T}}_{i,j} + \alpha\varepsilon - \hat{\mu}_i)^T + \beta\mathbf{E}, \quad (5)$$

259 where  $\varepsilon$  is learnable matrix initialized by randomly gaussian  
260 noise,  $\mathbf{E}$  is unit matrix,  $\alpha, \beta$  are hyper-parameters,  $\hat{\sigma}$  is tied  
261 covariance matrix, and  $\hat{\sigma} = [\hat{\sigma}_1, \hat{\sigma}_2, \dots, \hat{\sigma}_K]^T$ .

262 After computing the empirical Gaussian mean  $\hat{\mu}$  and  
263 the tied covariance matrix  $\hat{\sigma}$ , the Prompt Learning Module  
264 samples virtual OOD prompts from the low-likelihood re-  
265 gions of the class-conditional Gaussian distributions in hid-  
266 den space, based on the estimated multivariate distributions.  
267 Then, it selects the top-k prompts with the lowest probabili-  
268 ty from this  $\varepsilon$ -likelihood region:  
269

$$270 \quad \mathcal{V}_i = \Psi(\hat{\mathbf{T}}, \hat{\mu}_i, \hat{\sigma}), \quad (6)$$

271 where  $\Psi$  is class-conditional gaussian distribution prob-  
272 ability density and satisfies the following relation:

$$\begin{aligned} 273 \quad \Psi(\hat{\mathbf{T}}, \hat{\mu}_1, \hat{\mu}_2, \dots, \hat{\mu}_K, \hat{\sigma}) = \\ 274 \quad \Psi(\hat{\mathbf{T}}, \hat{\mu}_1, \hat{\sigma}) \Psi(\hat{\mathbf{T}}, \hat{\mu}_2, \hat{\sigma}) \dots \Psi(\hat{\mathbf{T}}, \hat{\mu}_K, \hat{\sigma}), \end{aligned} \quad (7)$$

275 For each  $\Psi(\hat{\mathbf{T}}, \hat{\mu}_i, \hat{\sigma})$ , its expansion can be formulated  
276 as:

$$277 \quad \Psi(\hat{\mathbf{T}}, \hat{\mu}_i, \hat{\sigma}) = \{v_i | \frac{1}{\sqrt{2\pi\hat{\sigma}}} e^{-\frac{1}{2}(v_i - \hat{\mu}_i)^T \hat{\sigma}^{-1} (v_i - \hat{\mu}_i)} < \epsilon\}, \quad (8)$$

278 where  $v_i \sim \mathcal{N}(\hat{\mu}_i, \hat{\sigma})$  denotes sampled virtual prompt  
279 using  $i_{th}$  ID category prompts,  $i = \{1, 2, \dots, K\}$ , and “ $^{-1}$ ”  
280 denotes matrix inverse operation. The final synthesised  
281 OOD prompts are denoted as  $\hat{\mathbf{T}}^\dagger$ .

### 282 3.2. OOD Virtual Images Synthesis

283 Existing methods [12, 13, 33–35] directly extract OOD  
284 pseudo-samples from ID data. However, the extracted  
285 pseudo-samples are unable to fit the distribution of OOD  
286 data adequately. In this paper, we also use synthesis method  
287 to get OOD data. The principle of synthesizing OOD im-  
288 age is similar to Eq. 3 to Eq. 8. Compared with synthesiz-  
289 ing OOD prompts, the input for calculating the empirical  
290 Gaussian mean and tied covariance is ID image embedding  
291 instead of ID prompts embedding. We define the final syn-  
292 thesised virtual images using current ID image embedding  
293 queue  $\mathcal{Q}_I$  as  $\hat{\mathbf{X}}^\dagger$ .

### 294 3.3. Text-Image Alignment Module

295 We first encode ID and OOD images and prompts to gen-  
296 erate their embeddings. Then, the similarity score between

296 prompts embedding and image embeddings is computed as  
297 follows:

$$298 \quad \mathbf{S} = \frac{\|\hat{\mathbf{X}}\|_p(\|\hat{\mathbf{T}}\|_p)^T}{e^\omega}, \quad (9)$$

299 where  $\hat{\mathbf{X}}$  is embedding of detected region images in the sec-  
300 ond training phase, and  $\hat{\mathbf{X}} = [\hat{\mathbf{X}}_1, \hat{\mathbf{X}}_2, \dots, \hat{\mathbf{X}}_b]$ ,  $\hat{\mathbf{X}}_i$  is one  
301 detected region image embedding,  $\hat{\mathbf{X}}_i \in \mathbb{R}^{l'}$ . In the third  
302 training phase, the input is embedding of synthesised vir-  
303 tual image  $\hat{\mathbf{X}}^\dagger$  instead of  $\hat{\mathbf{X}}$ ,  $\omega$  is hyper-parameters for scal-  
304 ing.  $\mathbf{S}$  is similarity score. The prompts embedding in Eq. 9  
305 is ID prompts embedding  $\hat{\mathbf{T}}$  in the second phase and synthe-  
306 sised OOD prompts embedding  $\hat{\mathbf{T}}^\dagger$  in the third phase,  $\|\cdot\|_p$   
307 is normalization, in addition,  $\|\hat{\mathbf{X}}_i\|_p = \hat{\mathbf{X}}_i / \sqrt{\sum_{j=1}^{l'} |\hat{\mathbf{X}}_{i,j}|^2}$   
308 and  $\|\hat{\mathbf{T}}_i^\dagger\|_p = \hat{\mathbf{T}}_i^\dagger / \sqrt{\sum_{j=1}^{l'} |\hat{\mathbf{T}}_{i,j}^\dagger|^2}$ .

### 309 3.4. Loss Function

310 Alignment loss  $\mathcal{L}_{align}$  constrains the contrastive learning  
311 process during alignment, receiving ID or OOD data at  
312 different training phases. The similarity score between  
313 prompts embedding and image embeddings is used to cal-  
314 culate the alignment loss:

$$315 \quad \mathcal{L}_{align}(\mathbf{S}, y) = - \sum_{i=1}^{K'} t_i \log(\mathcal{R}_i(\mathbf{S})), \quad (10)$$

316 where  $t_i$  represents category label of the object contained  
317 in currently detected region.  $\mathcal{R}_i$  represents the standardized  
318 prediction score. We treat all OOD categories as a single  
319 category, i.e., “background”. During the training phase, if  
320 the ID dataset contains a total of  $K$  classes, each detected  
321 region image is required to calculate similarity scores with  
322  $(K + 1)$  text prompts, i.e.,  $K' = K + 1$ .

323 Previous methods typically generate simple prompts  
324 that lack location information, such as “a photo of a  
325 <CLS>” [48, 49], or provide brief prompts with relative lo-  
326 cation information for the entire image [42]. We argue that  
327 these prompts lack the fine granularity needed for the model  
328 to learn essential location information in vision-language-  
329 based detection tasks.  $\mathcal{L}_{loc}$  is designed to implicitly incor-  
330 porate location information, enabling the generation of fine-  
331 grained prompts for detected image regions.

$$332 \quad \mathcal{L}_{loc} = \frac{\lambda}{\Phi(\mathbf{B}_g)} \left[ \sum_{i=1}^z (\sqrt{\mathbf{B}_{g_i}} - \sqrt{u(\mathbf{B}_r)_i})^2 \right]^{\frac{1}{2}}, \quad (11)$$

333 where  $\mathbf{B}_g$  represents ground truth coordinates of detected  
334 image region,  $\mathbf{B}_r$  represents regression results of coordi-  
335 nates, and  $\mathbf{B}_g \in \mathbb{R}^{b \times 4}$ ,  $\mathbf{B}_r \in \mathbb{R}^{b \times 4}$ ,  $z = 4$ ,  $u$  represents  
336 calculating absolute values,  $\Phi$  represents calculating the di-  
337 mension of vector,  $\lambda$  is hyper-parameter.

After incorporating the classification loss  $\mathcal{L}_{cls}$  and the  
location loss  $\mathcal{L}_{loc}$ , the total loss can be expressed as:

$$\begin{aligned} \mathcal{L} = & \xi_1 [\gamma_1 \tau \mathcal{L}_{align}^{id} + \gamma_2 (1 - \tau) \mathcal{L}_{align}^{ood}] \\ & + \gamma_3 \xi_2 [\kappa \mathcal{L}_{loc}^{id} + (1 - \kappa) \mathcal{L}_{loc}^{ood}] \\ & + \gamma_4 \xi_3 \mathcal{L}_{cls} + \gamma_5 \xi_4 \mathcal{L}_{reg} + \overline{\mathcal{W}}. \end{aligned} \quad (12)$$

Note that  $\gamma_1, \gamma_2, \gamma_3, \gamma_4, \gamma_5$  are the hyper-parameters,  $\xi, \tau, \kappa$   
determine the loss functions used in the current training  
phase and  $\xi_i = \{0, 1\}$ ,  $\tau = \{0, 1\}$ . In order to bet-  
ter regularize the model, in the actual implementation of  
 $\mathcal{L}$ , we also add the regularization term  $\mathcal{W}$ , and  $\mathcal{W}_i =$   
 $[\Delta_{(\mathcal{F}(\mathbf{O}_1), \mathcal{B}_1)_i}^2 + \Delta_{(\mathcal{G}(\mathbf{O}_2), \mathcal{B}_2)_i}^2]$ ,  $\overline{\mathcal{W}} = 1/N \sum_{i=1}^N \mathcal{W}_i$ ,  $\Delta_{(a,b)}^2$   
represents  $(a - b)^2$ ,  $\mathcal{F}, \mathcal{G}$  represent regression blocks,  $\mathbf{O}_i$   
represents regularization matrix,  $\mathcal{B}_i$  represents bias matrix  
of regression block,  $i = \{1, 2\}$ .

## 4. Experiments

### 4.1. Datasets

We verify our proposed APLGOS on four commonly used  
datasets: PASCAL VOC, Berkley DeepDrive-100k, MS-  
COCO2017 and OpenImages. The PASCAL VOC [7]  
dataset contains 9963 images in 20 categories, split into  
5011 training and 4952 test images, with a resolution of  
 $500 \times 375$  ( $375 \times 500$ ). The BDD-100k [43] dataset con-  
sists of 100,000 high-resolution driving scenarios with de-  
tailed road object annotations. The MS-COCO2017 [24]  
dataset includes 328,000 images across 91 categories and  
2.5 million instance tags, with 82 categories having more  
than 5000 tags. OpenImages V4 [15] contains 9.2 million  
images across 500 categories, commonly used for classifi-  
cation, object detection, and visual relationship detection.  
The above four datasets comprehensively evaluate our pro-  
posed method from different aspects and perspectives.

### 4.2. Implementation Details

We employ transformer as the backbone for the text en-  
coder in the Prompt Learning Module. For the image en-  
coder, we employ ResNet50 [10] and RegNetX4.0 [31] as  
backbones, respectively. We use ChatGPT-3.5 to standard-  
ize Q&A pairs. The ratio of ID data used for training to  
synthesised OOD data is approximately 1:1. We use PAS-  
CAL VOC and Berkeley DeepDrive-100K as ID datasets,  
and evaluate on two OOD datasets containing subsets ran-  
domly sampled from MS-COCO2017 and OpenImages, re-  
spectively. To ensure the fairness of the test, we manu-  
ally exclude the categories in the OOD dataset that overlap  
with those in the ID dataset before evaluating on the OOD  
dataset. We set  $B = 16$  and train APLGOS on PASCAL  
VOC for 18,000 iterations, and set  $B = 8$  to train on Berke-  
ley DeepDrive-100k for 90,000 iterations. We set the learn-  
ing rate  $lr = 0.01$ . The length of prompt embedding and

ID Dataset	Method	FPR95 ↓	AUROC ↑	AUPR ↑	mAP (ID) ↑
			OOD: MS-COCO2017 / OpenImages		
PASCAL VOC	MSP [11]	70.99 / 73.13	83.45 / 81.91	-	48.7
	ODIN [22]	59.82 / 63.14	82.20 / 82.59	-	48.7
	Mahalanobis [19]	96.46 / 96.27	59.25 / 57.42	-	48.7
	Energy score [25]	56.89 / 58.69	83.69 / 82.98	-	48.7
	Gram matrices [32]	62.75 / 67.42	79.88 / 77.62	-	48.7
	Generalized ODIN [14]	59.57 / 70.28	83.12 / 79.23	-	48.1
	CSI [35]	59.91 / 57.41	81.83 / 82.95	-	48.1
	GAN-synthesis [18]	60.93 / 59.97	83.67 / 82.67	-	48.5
	VOS-ResNet50* [6]	48.28 / 52.14	87.65 / 85.3	98.76 / 96.98	47.8
	VOS-RegX4.0* [6]	50.53 / 50.27	88.10 / 87.08	98.92 / 97.80	49.1
	<b>APLGOS (ResNet50)</b>	<b>47.16 / 49.66</b>	<b>87.89 / 85.91</b>	<b>98.80 / 97.54</b>	<b>48.8</b>
<b>APLGOS (RegNetX4.0)</b>	<b>45.96 / 47.10</b>	<b>89.19 / 88.49</b>	<b>99.00 / 98.30</b>	<b>49.4</b>	
Berkeley DeepDrive-100k	MSP [11]	80.94 / 79.04	75.87 / 77.38	-	31.2
	ODIN [22]	62.85 / 58.92	74.44 / 76.61	-	31.2
	Mahalanobis [19]	57.66 / 60.16	84.92 / 86.88	-	31.2
	Energy score [25]	60.06 / 54.97	77.48 / 79.60	-	31.2
	Gram matrices [32]	60.93 / 77.55	74.93 / 59.38	-	31.2
	Generalized ODIN [14]	57.27 / 50.17	85.22 / 87.18	-	31.8
	CSI [35]	47.10 / 37.06	84.09 / 87.99	-	30.6
	GAN-synthesis [18]	57.03 / 50.61	78.82 / 81.25	-	31.4
	VOS-ResNet50* [6]	46.97 / 31.25	84.97 / 89.82	99.67 / 99.86	35.7
	VOS-RegX4.0* [6]	42.82 / 27.55	86.36 / 92.11	99.76 / 99.93	37.0
	Dynamic Prototypes [37]	45.72 / 35.05	85.14 / 88.92	-	31.5
	<b>APLGOS (ResNet50)</b>	<b>41.10 / 23.30</b>	<b>87.36 / 92.87</b>	<b>99.73 / 99.89</b>	<b>35.8</b>
	<b>APLGOS (RegNetX4.0)</b>	<b>39.48 / 19.79</b>	<b>87.47 / 93.59</b>	<b>99.79 / 99.94</b>	<b>37.6</b>

Table 1. Comparison with the state-of-the-art methods on mainstream datasets. Here we use PASCAL VOC and Berkeley DeepDrive-100k as ID datasets, MS-COCO2017 and OpenImages as OOD datasets, respectively. “-” denotes that the data is not available.

Strategy	FPR95 ↓	AUROC ↑	AUPR ↑	mAP (ID) ↑
		OOD: MS-COCO2017 / OpenImages		
(a) VOS-RegNetX4.0* [6]	50.53 / 50.27	88.10 / 87.08	98.82 / 97.80	49.1
(b) [6] + <CLS>	50.12 / 49.50	88.56 / 86.83	98.91 / 97.79	48.2
(c) [6] + “a region of a” + <CLS>	51.31 / 50.96	88.20 / 86.73	98.98 / 97.85	48.7
(d) [6] + RP + <CLS>	49.50 / 49.40	88.49 / 86.73	98.82 / 97.77	48.9
(e) [6] + <LOC> + “a region of a” + <CLS>	49.56 / 47.60	88.23 / 87.07	98.89 / 97.87	49.1
(f) [6] + <LOC> + RP + <CLS> (Ours)	<b>45.96 / 47.10</b>	<b>89.19 / 88.49</b>	<b>99.00 / 98.30</b>	<b>49.4</b>

Table 2. Ablation studies for prompt strategies. “+” denotes the combination of strategies. “RP” represents sampled prompts from statements set, which is standardized by ChatGPT using Q&A pair and guidance instructions. (b) denotes the simplest prompt strategy, i.e., only providing the ground-truth label for the ID data, (for synthesised OOD image, we define its label as “background”). (c) denotes the original prompt strategy of CLIP [30]. (d) denotes that we replace the prompts in CLIP [30] with the statements by ChatGPT standardizing the Q&A pairs. (e) denotes adding location tokens <LOC> to (c). (f) represents the prompts of our proposed APLGOS.

length of image embedding  $l' = 1024$ . We use 1000 samples to estimate the class-conditional Gaussian distribution of ID image embeddings and 10000 samples for ID prompts embedding (i.e.,  $|Q_I| = 1000$ ,  $|Q_T| = 10000$ ). The total length  $l$  of the standardized Q&A pair does not exceed 77. In the experimental tables, “\*” denotes results from local replication based on open-source code. “↓” indicates that a smaller value is better, while “↑” indicates that a greater value is better.

### 4.3. Comparison with The State-of-the-Art

We report the results of our proposed framework with different image encoder backbones (ResNet50 and Reg-

NetX4.0) on PASCAL VOC, Berkeley DeepDrive-100k, MS-COCO2017, and OpenImages datasets, as shown in Table 1. The best results for the same dataset and the same backbone settings are shown in **bold**. For the same evaluation metric on the same dataset, the best results are underlined. When using Transformer-based text encoder and ResNet50-based image encoder, APLGOS achieves an FPR95 of 47.16% and an mAP of 48.8% on PASCAL VOC (ID) with MS-COCO2017 as the OOD dataset. When OpenImages is used as the OOD dataset, FPR95 increases to 49.66%. Compared to the state-of-the-art OOD detection model [6], APLGOS reduces FPR95 by 1.12% and 2.48% on MS-COCO2017 and OpenImages, respectively.

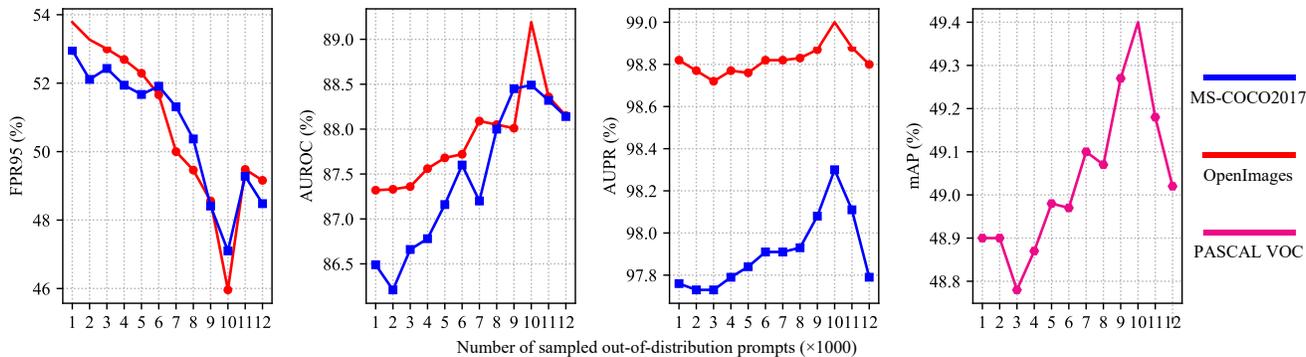


Figure 3. Ablation on number of sampled OOD prompts  $\mathcal{K}$ . The horizontal coordinate is the number of sampled ood prompts  $\mathcal{K} (\times 10^3)$ , while the vertical coordinates are, from left to right, FPR95, AUROC, AUPR, and mAP, respectively. **Red line** and **Blue line** represent using MS-COCO2017 and OpenImages as OOD datasets, respectively. **Pink line** represents using PASCAL VOC as ID dataset.

$\alpha$	FPR95 ↓	AUROC ↑	mAP (ID) ↑
	OOD: MS-COCO2017 / OpenImages		
0	51.63 / 50.88	87.86 / 87.24	49.2
0.5	51.90 / 51.48	87.55 / 87.02	48.9
<b>1.0</b>	<b>45.96 / 47.10</b>	<b>89.19 / 88.49</b>	<b>49.4</b>
1.5	55.88 / 53.33	86.29 / 86.75	48.9
2.0	55.92 / 49.54	86.75 / 88.00	48.9

Table 3. The Ablation Experiments on The Strength of Random Gaussian Noise  $\epsilon$ .  $\alpha$  represents the strength of added gaussian noise. The value of  $\alpha$  increases gradually from 0 to 2.0, and we take the value at 0.5 intervals.

$\Gamma_1$	FPR95 ↓	AUROC ↑	mAP (ID) ↑
	OOD: MS-COCO2017 / OpenImages		
1:4	50.11 / 58.38	87.71 / 85.67	49.1
1:3	49.40 / 55.12	87.91 / 86.38	49.2
1:2	47.98 / 54.49	88.40 / 85.94	49.2
<b>1:1</b>	<b>45.96 / 47.10</b>	<b>89.19 / 88.49</b>	<b>49.4</b>
2:1	48.25 / 50.20	88.30 / 87.76	49.2
3:1	50.95 / 53.94	86.81 / 84.70	47.5
4:1	50.20 / 51.56	86.70 / 84.89	47.3

Table 4. The ablation experiments on the ratio  $\Gamma_1$  of ID and OOD data used during training. Our default parameters and results are shown in **bold**. Parameters and results of baseline [6] are shown with a dark base color.

409 With Transformer-based text encoder and RegNetX4.0-  
 410 based image encoder, FPR95 decreases to 45.96% on MS-  
 411 COCO2017 and 47.1% on OpenImages, while the mAP  
 412 on PASCAL VOC improves to 49.4%. This setup further  
 413 reduces FPR95 by 4.57% and 3.17% on MS-COCO2017  
 414 and OpenImages, respectively, compared to [6]. For  
 415 Berkley DeepDrive-100k (ID), using ResNet50-based im-  
 416 age encoder and Transformer-based text encoder, APL-  
 417 GOS achieves an FPR95 of 41.10% on MS-COCO2017 and  
 418 23.30% on OpenImages, with an mAP of 35.8%. When us-  
 419 ing RegNetX4.0-based image encoder instead, FPR95 fur-  
 420 ther decreases to 39.48% on MS-COCO2017 and 19.79%  
 421 on OpenImages, while mAP improves to 37.6%.

#### 422 4.4. Ablation Studies

423 **Prompt strategies.** To further validate the effectiveness of  
 424 our prompt strategies, we conduct extensive ablation exper-  
 425 iments on APLGOS’s prompt strategies, and the results are  
 426 shown in Table 2. Sampling from the statements set brings  
 427 greater performance gains than simply initializing learnable  
 428 prompts with “a region of a” ((c) vs (d)). Moreover, adding  
 429 location tokens to prompts significantly improves perfor-  
 430 mance, as it refines the scope of the prompts ((c) vs (e)).  
 431 Compared to other prompt strategies, our APLGOS prompt  
 432 strategy (f) integrates the advantages of the aforementioned

strategies and achieves the best performance. 433

**Number of Sampled OOD Prompts.** APLGOS synthe- 434  
 sises virtual prompts for OOD categories and for each ID 435  
 category, APLGOS samples  $\mathcal{K}$  virtual OOD prompts in low- 436  
 likelihood regions of ID class-conditional gaussian distribu- 437  
 tions in high-dimensional hidden space. We conduct ablation 438  
 experiments on  $\mathcal{K}$ , the results of its effect on perfor- 439  
 mance are shown in Figure 3. When  $\mathcal{K}$  is too small, it may 440  
 fail to adequately cover the region outside the ID categories’ 441  
 decision boundaries in the hidden space. On the other hand, 442  
 when  $\mathcal{K}$  is too large, the excessive randomness in the sam- 443  
 pled OOD prompts makes it difficult to effectively regular- 444  
 ize the decision boundaries with the limited model param- 445  
 eters. Therefore, we set  $\mathcal{K} = 10000$  as the default value. 446

**Strength of Random Gaussian Noise  $\epsilon$ .** To enhance the 447  
 size and diversity of the OOD prompts embedding sampling 448  
 space and prevent the model from overly relying on the 449  
 ID category distribution, we introduce a learnable matrix 450  
 initialized with random Gaussian noise  $\epsilon$  during the OOD 451  
 prompt sampling stage (Eq. 5). We conduct ablation exper- 452  
 iments on its strength  $\alpha$ , and the results are shown in 453  
 Table 3. A small value of  $\alpha$  makes the sampling space of 454  
 OOD prompts embedding too narrow, while a large value 455

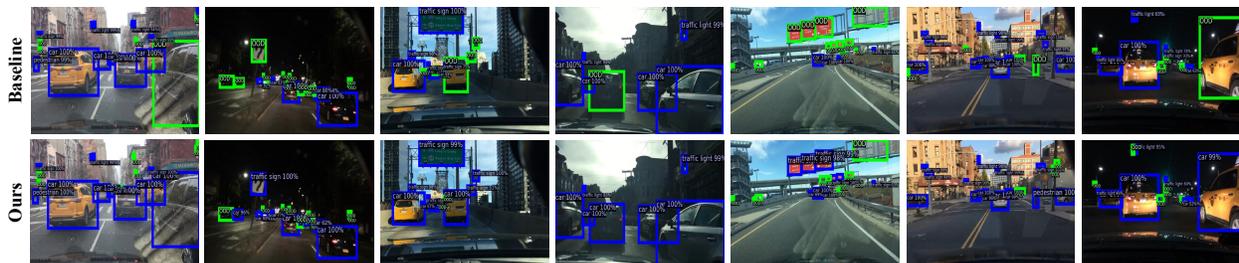


Figure 4. Detection results on ID dataset. Here we use Berkley DeepDrive-100k dataset as ID dataset. We use RegNetX4.0 and Transformer as backbone. The **first row** is the detection results of baseline [6]. The **second row** is the detection result of our APLGOS. Our APLGOS rarely misclassifies the ID class as OOD class, and there is almost no missed detection.

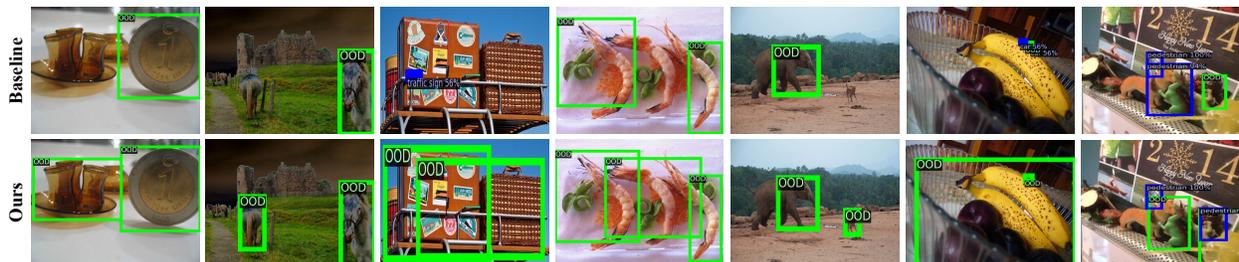


Figure 5. Detection results on OOD datasets. Here we use Berkley DeepDrive-100k dataset as ID dataset, MS-COCO2017 and OpenImages as OOD datasets. The **first row** is the detection results of baseline [6]. The **second row** is the detection results of our APLGOS. Compared to the baseline, APLGOS rarely misses detections and hardly produces overlapping boxes for the same object.

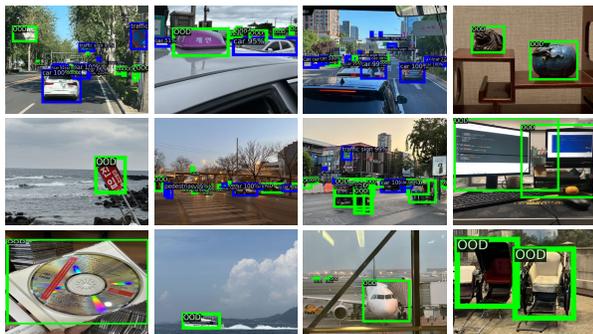


Figure 6. Detection results in Real World. Here we use Berkley DeepDrive-100k dataset as in-distribution dataset. Pictures we take ourselves with our phone as out-of-distribution dataset.

456 of  $\alpha$  results in an overly large sampling space. Only by ap-  
 457 propriately expanding the sampling space of OOD prompts  
 458 embedding can the model’s ability to fit the OOD distribu-  
 459 tion be effectively enhanced.

460 **Ratio of ID and OOD Data Used During Training.**

461 To verify that APLGOS can achieve better performance  
 462 with less ID data, we conduct ablation experiments on the  
 463 amount of ID data used during training, and the results are  
 464 shown in Table 4. By default, APLGOS adopts a ratio  $\Gamma_1$   
 465 of 1:1 for ID and OOD data during training, whereas the  
 466 baseline [6] uses a ratio of 2:1. However, in this case, the  
 467 performance of APLGOS decreases instead.

468 **Visualization of Detection Results.** To better evaluate the  
 469 performance of APLGOS, we visualize its detection results  
 470 on ID datasets, OOD datasets, and real-world scenarios.

The results are presented in Figures 4, 5 and 6. The images  
 471 in real-world scenarios are captured using an iPhone 14 Pro  
 472 Max. The visualization results demonstrate that APLGOS  
 473 outperforms the baseline method in detecting ID and OOD  
 474 categories. Moreover, the visualization of detection results  
 475 in real-world scenarios further confirms its superior gener-  
 476 alization ability.  
 477

478 **5. Conclusion**

In this paper, we propose a vision-language method, Adap-  
 479 tive Prompt Learning via Gaussian Outlier Synthesis (APL-  
 480 GOS) for Out-of-distribution Detection. Through prompt  
 481 learning approach, APLGOS provides adaptive region-level  
 482 prompts with location information for ID / OOD images.  
 483 We use ChatGPT to standardize pre-defined Q&A pairs and  
 484 generate a statements set. During training, only ID im-  
 485 ages are from the dataset, while ID prompts, OOD prompts  
 486 and OOD images are all virtual. We sample statements  
 487 from the statements set to initialize learnable ID prompts.  
 488 We samples virtual OOD prompts and OOD images in  
 489 the low-likelihood region of the class-conditional gaussian  
 490 distribution in high-dimensional hidden space. Similarity  
 491 score between prompts and images is utilized to calculate  
 492 contrastive learning loss in high-dimensional hidden space,  
 493 which guarantees the quality of virtual outliers as well as  
 494 better regularization of the model. Through comprehen-  
 495 sive experimental evaluations, we demonstrated the effec-  
 496 tiveness of the proposed APLGOS.  
 497

498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553

**References**

[1] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. Aleman, D. Almeida, J. Altschmidt, S. Altman, S. Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 2

[2] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems (NeurIPS)*, 33:1877–1901, 2020. 2

[3] L. Chen, G. Wang, L. Yuan, K. Wang, K. Deng, and P. Torr. Nerf-vpt: Learning novel view representations with neural radiance fields via view prompt tuning. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pages 1156–1164, 2024. 2

[4] J. Devlin, M. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 2

[5] X. Du, Z. Wang, M. Cai, and S. Li. Towards unknown-aware learning with virtual outlier synthesis. In *International Conference on Learning Representations (ICLR)*, 2022. 1

[6] X. Du, Z. Wang, M. Cai, and Y. Li. Vos: Learning what you don’t know by virtual outlier synthesis. *International Conference on Learning Representations (ICLR)*, 2022. 2, 6, 7, 8

[7] M. Everingham, L. Van Gool, C. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision (IJCV)*, 88:303–338, 2010. 5

[8] Zhongbin Fang, Xiangtai Li, Xia Li, Joachim M Buhmann, Chen Change Loy, and Mengyuan Liu. Explore in-context learning for 3d point cloud understanding. *Advances in Neural Information Processing Systems*, 36, 2024. 2

[9] M. Grcić, P. Bevandić, and S. Šegvić. Dense open-set recognition with synthetic outliers generated by real nvp. *arXiv preprint arXiv:2011.11094*, 2020. 2

[10] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 5

[11] D. Hendrycks and K. Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *International Conference on Learning Representations (ICLR)*, 2016. 6

[12] D. Hendrycks, M. Mazeika, and T. Dietterich. Deep anomaly detection with outlier exposure. *International Conference on Learning Representations (ICLR)*, 2018. 1, 4

[13] D. Hendrycks, S. Basart, M. Mazeika, A. Zou, J. Kwon, M. Mostajabi, J. Steinhardt, and D. Song. Scaling out-of-distribution detection for real-world settings. *arXiv preprint arXiv:1911.11132*, 2019. 1, 2, 4

[14] Y. Hsu, Y. Shen, H. Jin, and Z. Kira. Generalized odin: Detecting out-of-distribution image without learning from out-of-distribution data. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)*, pages 10951–10960, 2020. 2, 6

[15] A. Kuznetsova, H. Rom, N. Alldrin, J. Uijlings, I. Krasin, J. Pont-Tuset, S. Kamali, S. Popov, M. Mallocci, A. Kolesnikov, et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *International Journal of Computer Vision (IJCV)*, 128(7):1956–1981, 2020. 5

[16] T. Le Scao, A. Fan, C. Akiki, E. Pavlick, S. Ilić, D. Hesslow, R. Castagné, A. Luccioni, F. Yvon, M. Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. 2023. 2

[17] Y. LeCun, S. Chopra, R. Hadsell, M. Ranzato, and F. Huang. A tutorial on energy-based learning. *Predicting structured data*, 1(0), 2006. 2

[18] K. Lee, H. Lee, K. Lee, and J. Shin. Training confidence-calibrated classifiers for detecting out-of-distribution samples. *International Conference on Learning Representations (ICLR)*, 2017. 2, 6

[19] K. Lee, K. Lee, H. Lee, and J. Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *Advances in neural information processing systems (NeurIPS)*, 31, 2018. 6

[20] D. Li, J. Li, and S. Hoi. Blip-diffusion: Pre-trained subject representation for controllable text-to-image generation and editing. *Advances in Neural Information Processing Systems (NeurIPS)*, 36, 2024. 1

[21] J. Li, P. Chen, Z. He, S. Yu, S. Liu, and J. Jia. Rethinking out-of-distribution (ood) detection: Masked image modeling is all you need. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)*, pages 11578–11589, 2023. 1

[22] S. Liang, Y. Li, and R. Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. *International Conference on Learning Representations (ICLR)*, 2017. 2, 6

[23] J. Liao, X. Xu, M. Nguyen, A. Goodge, and C. Foo. Coftad: Contrastive fine-tuning for few-shot anomaly detection. *IEEE Transactions on Image Processing (TIP)*, 2024. 1

[24] T. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision (ECCV)*, pages 740–755. Springer, 2014. 5

[25] W. Liu, X. Wang, J. Owens, and Y. Li. Energy-based out-of-distribution detection. *Advances in neural information processing systems (NeurIPS)*, 33:21464–21475, 2020. 2, 6

[26] X. Liu, Y. Lochman, and C. 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 (CVPR)*, pages 23946–23955, 2023. 2

[27] Yifei Ming, Ziyang Cai, Jiuxiang Gu, Yiyun Sun, Wei Li, and Yixuan Li. Delving into out-of-distribution detection with vision-language representations. *Advances in Neural Information Processing Systems (NeurIPS)*, 35:35087–35102, 2022. 2

[28] S. Mohseni, M. Pitale, J. Yadawa, and Z. Wang. Self-supervised learning for generalizable out-of-distribution detection. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pages 5216–5223, 2020. 2

554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610

611 [29] S. Ning, L. Qiu, Y. Liu, and X. He. Hoiclip: Efficient knowl- 669  
612 edge transfer for hoi detection with vision-language mod- 670  
613 els. In *Proceedings of the IEEE/CVF Conference on Com- 671*  
614 *puter Vision and Pattern Recognition (CVPR)*, pages 23507– 672  
615 23517, 2023. 2 673

616 [30] A. Radford, J. Kim, C. Hallacy, A. Ramesh, G. Goh, S. 674  
617 Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. 675  
618 Learning transferable visual models from natural language 676  
619 supervision. In *International conference on machine learn- 677*  
620 *ing (ICML)*, pages 8748–8763, 2021. 2, 6 678

621 [31] I. Radosavovic, R. Kosaraju, R. Girshick, K. He, and P. 679  
622 Dollár. Designing network design spaces. In *Proceedings 680*  
623 *of the IEEE/CVF Conference on Computer Vision and Pat- 681*  
624 *tern Recognition (CVPR)*, pages 10428–10436, 2020. 5 682

625 [32] C. Sastry and S. Oore. Detecting out-of-distribution exam- 683  
626 ples with gram matrices. In *International Conference on Ma- 684*  
627 *chine Learning (ICML)*, pages 8491–8501, 2020. 1, 6 685

628 [33] B. Su, H. Zhang, and Z. Zhou. Hsic-based moving weight 686  
629 averaging for few-shot open-set object detection. In *Pro- 687*  
630 *ceedings of the 31st ACM International Conference on Mul- 688*  
631 *timedia (ACM MM)*, pages 5358–5369, 2023. 1, 4 689

632 [34] B. Su, H. Zhang, J. Li, and Z. Zhou. Toward generalized 690  
633 few-shot open-set object detection. *IEEE Transactions on 691*  
634 *Image Processing (TIP)*, 33:1389–1402, 2024. 692

635 [35] J. Tack, S. Mo, J. Jeong, and J. Shin. Csi: Novelty detec- 693  
636 tion via contrastive learning on distributionally shifted in- 694  
637 stances. *Advances in neural information processing systems 695*  
638 *(NeurIPS)*, 33:11839–11852, 2020. 1, 2, 4, 6 696

639 [36] H. Wang, Y. Li, H. Yao, and X. Li. Clipn for zero-shot 697  
640 ood detection: Teaching clip to say no. In *Proceedings of 698*  
641 *the IEEE/CVF International Conference on Computer Vision 699*  
642 *(ICCV)*, pages 1802–1812, 2023. 2 700

643 [37] A. Wu, C. Deng, and W. Liu. Unsupervised out-of- 701  
644 distribution object detection via pca-driven dynamic proto- 702  
645 type enhancement. *IEEE Transactions on Image Processing 703*  
646 *(TIP)*, 2024. 6 704

647 [38] X. Wu, F. Zhu, R. Zhao, and H. Li. Cora: Adapting clip 705  
648 for open-vocabulary detection with region prompting and an- 706  
649 chor pre-matching. In *Proceedings of the IEEE/CVF confer- 707*  
650 *ence on computer vision and pattern recognition (CVPR)*, 708  
651 pages 7031–7040, 2023. 2 709

652 [39] C. Xie, Z. Zhang, Y. Wu, F. Zhu, R. Zhao, and S. Liang. 710  
653 Described object detection: Liberating object detection with 711  
654 flexible expressions. *Advances in Neural Information Pro- 712*  
655 *cessing Systems (NeurIPS)*, 36, 2024. 1 713

656 [40] X. Xiong, Y. Liu, T. Yuan, Y. Wang, Y. Wang, and H. Zhao. 714  
657 Neural map prior for autonomous driving. In *Proceedings of 715*  
658 *the IEEE/CVF Conference on Computer Vision and Pattern 716*  
659 *Recognition (CVPR)*, pages 17535–17544, 2023. 1 717

660 [41] J. Xu, X. Liu, Y. Wu, Y. Tong, Q. Li, M. Ding, J. Tang, and Y. 718  
661 Dong. Imagereward: Learning and evaluating human pref- 719  
662 erences for text-to-image generation. *Advances in Neural 720*  
663 *Information Processing Systems (NeurIPS)*, 36, 2024. 1 721

664 [42] Z. Yang, Z. Gan, J. Wang, X. Hu, F. Ahmed, Z. Liu, Y. Lu, 722  
665 and L. Wang. Unitab: Unifying text and box outputs for 723  
666 grounded vision-language modeling. In *European Confer- 724*  
667 *ence on Computer Vision (ECCV)*, pages 521–539. Springer, 725  
668 2022. 5 726

[43] F. Yu, H. Chen, X. Wang, W. Xian, Y. Chen, F. Liu, V. Mad- 669  
havan, and T. Darrell. Bdd100k: A diverse driving dataset 670  
for heterogeneous multitask learning. In *Proceedings of 671*  
*the IEEE/CVF conference on computer vision and pattern 672*  
*recognition (CVPR)*, pages 2636–2645, 2020. 5 673

[44] Y. Zhai, Y. Zeng, Z. Huang, Z. Qin, X. Jin, and D. 674  
Cao. Multi-prompts learning with cross-modal alignment for 675  
attribute-based person re-identification. In *Proceedings of 676*  
*the AAAI Conference on Artificial Intelligence (AAAI)*, pages 677  
6979–6987, 2024. 2 678

[45] C. Zhang, X. Chen, S. Chai, C. Wu, D. Lagun, T. Beeler, and 679  
F. De la Torre. Iti-gen: Inclusive text-to-image generation. 680  
In *Proceedings of the IEEE/CVF International Conference 681*  
*on Computer Vision (ICCV)*, pages 3969–3980, 2023. 2 682

[46] J. Zhang, N. Inkawhich, Y. Chen, and H. Li. Fine-grained 683  
out-of-distribution detection with mixup outlier exposure. 684  
*arXiv preprint arXiv:2106.03917*, 2(5), 2021. 2 685

[47] J. Zhang, J. Huang, S. Jin, and S. Lu. Vision-language mod- 686  
els for vision tasks: A survey. *IEEE Transactions on Pattern 687*  
*Analysis and Machine Intelligence (TPAMI)*, 2024. 2 688

[48] K. Zhou, J. Yang, C. Loy, and Z. Liu. Conditional prompt 689  
learning for vision-language models. In *Proceedings of 690*  
*the IEEE/CVF conference on computer vision and pattern 691*  
*recognition (CVPR)*, pages 16816–16825, 2022. 2, 5 692

[49] K. Zhou, J. Yang, C. Loy, and Z. Liu. Learning to prompt for 693  
vision-language models. *International Journal of Computer 694*  
*Vision (IJCV)*, 130(9):2337–2348, 2022. 2, 5 695

[50] Z. Zhu, Y. Zhang, H. Chen, Y. Dong, S. Zhao, W. Ding, J. 696  
Zhong, and S. Zheng. Understanding the robustness of 3d 697  
object detection with bird’s-eye-view representations in au- 698  
tonomous driving. In *Proceedings of the IEEE/CVF Confer- 699*  
*ence on Computer Vision and Pattern Recognition (CVPR)*, 700  
pages 21600–21610, 2023. 1 701

[51] O. Zohar, K. Wang, and S. Yeung. Prob: Probabilistic ob- 702  
jectness for open world object detection. In *Proceedings of 703*  
*the IEEE/CVF Conference on Computer Vision and Pattern 704*  
*Recognition (CVPR)*, pages 11444–11453, 2023. 1 705