

分 数:	
评卷人:	

华中科技大学

研究生“高级机器学习”课程论文(报告)

题 目: 基于 CLIP 的小样本分布外检测

学 号 U202115210

姓 名 梁一凡

专 业 人工智能本硕博 2101 班

课程指导教师 伍冬睿 朱力军

院 (系、所) 人工智能与自动化学院

2025 年 5 月 30 日

基于 CLIP 的小样本分布外检测

梁一凡

人工智能与自动化学院

摘要: 在开放世界应用中, 分布外 (OOD) 检测对于部署可靠的机器学习模型至关重要。常见的 OOD 检测方法使用单模态监督学习, 这些方法取得了不错的成果, 但存在局限性。如今以 CLIP 为代表的视觉语言模型在各种下游任务中取得了令人印象深刻的性能, 将 CLIP 应用于 OOD 检测也开始引起越来越多的关注, 因此本文提出了一种基于 CLIP 的小样本 OOD 检测方法。

小样本 OOD 检测旨在仅使用少量带标签的分布内 (ID) 图像, 检测出训练中未见过的 OOD 图像。尽管像 CoOp 这样的提示学习方法在少样本 ID 分类中展现了有效性, 但由于 OOD 背景信息干扰了这一学习过程, 其在 OOD 检测中仍存在局限性。为解决这一问题, 我们引入了一种名为 LoCoOp 的方法, 该方法在训练期间利用 OOD 背景的局部特征进行 OOD 正则化, 抑制了 OOD 背景的干扰, 进而增强 ID 与 OOD 之间的区分度。然而, 由于 ID 前景与 OOD 背景的分解不准确, 从 ID 数据中挖掘的 OOD 背景可能存在虚假性, 从而限制了 OOD 检测性能, 因此本文进一步提出名为 SCT 的训练框架, 对原始学习目标的两个部分引入调制因子, 在训练过程中根据预测的不确定性, 自适应地引导两个任务之间的优化过程, 以校准 OOD 正则化的影响。此外, 注意到当前检测范式没有充分利用训练集的视觉信息, 因此引入双模式匹配分数 DPM, 结合视觉匹配和文本匹配以期取得更好的检测效果。

在大规模 ImageNet OOD 检测基准上的实验表明, 本文的方法 LoCoOp 和 SCT 在小样本场景下均取得了较好的分布外检测性能, 当结合 DPM 分数时, CoOp、LoCoOp、SCT 方法均取得了进一步的性能提升, 尤其是在 iNaturalist 这个 OOD 数据集上。

代码和数据可以在 https://github.com/YifanLiang-hust/Advanced_ML_Final_Project 获得。纸质版的报告非最终版本, 请以电子版的报告为准。

关键词: 分布外检测; 对比式视觉语言模型 CLIP; 小样本学习

一、引言

在开放世界场景中，部署的机器学习模型可能会遇到训练期间未出现的分布外（OOD）样本。对 OOD 样本的错误预测可能会导致严重后果，尤其是在自动驾驶、智慧医疗等关键领域，因此检测 OOD 样本的能力至关重要。常见的 OOD 检测方法使用单模态监督学习，这些方法取得了不错的成果，但存在局限性，例如，这些方法在训练时需要大量的计算和标注成本。如今以 CLIP 为代表的视觉语言模型在各种下游任务中取得了令人印象深刻的性能，将 CLIP 应用于 OOD 检测也开始引起越来越多的关注。



图 1 机器学习模型可能会遇到训练期间未出现的分布外样本

以往基于 CLIP 的 OOD 检测的一些研究探索了无需任何分布内（ID）训练数据的零样本 OOD 检测方法，而其他研究则开发了需要完整的 ID 训练数据的完全监督方法。这些方法都取得了一定的成果，但两种方法都有各自的局限性。一方面，零样本方法不需要任何训练数据，但它们可能会与下游 ID 数据存在领域差距，这限制了零样本方法的性能。另一方面，完全监督方法利用整个 ID 训练数据，但通过微调可能会破坏 CLIP 的丰富表示，这也限制了其性能，尽管需要巨大的训练成本。为了克服完全监督和零样本方法的局限性，开发一种利用少量 ID 训练图像进行 OOD 检测的方法至关重要。

在本文中，我们专注于基于 CLIP 的小样本 OOD 检测问题，即仅使用少量的 ID 图像来进行 CLIP 的微调，其中最常见方法是提示学习，该方法在固定预训练参数的同时训练可学习的提示。CoOp 是提示学习的代表性工作，但其对 OOD 检测仍不太有效。虽然 CoOp 旨在提高给定的 ID 图像嵌入及其对应的类别文本嵌入的对齐程度，但它也可能使文本嵌入更接近 ID 图像中与 ID 无关的 OOD 背景。因此，这些文本嵌入可能包含与 ID 图像无关的信息，从而导致对 OOD 图像的置信度得分出现错误的过高情况。

为了将提示学习应用于 OOD 检测，我们提出了一种简单但有效的方法，称为

LoCoOp。我们观察到 CLIP 的局部特征包含许多与类别无关的干扰因素，如图 2 所示。LoCoOp 将这些与类别无关的干扰因素视为 OOD，并学习将它们从类别文本嵌入中推开，从而从类别文本嵌入中去除不必要的信息，并防止模型对异常特征生成较高的类别置信度分数。为了实现这一点，LoCoOp 首先使用每个局部区域的分割分类预测结果来识别与类别无关的区域。接下来，它对提取的区域的分类概率进行熵最大化，这确保了 ID 无关区域中的特征与任何类别文本嵌入都不同。

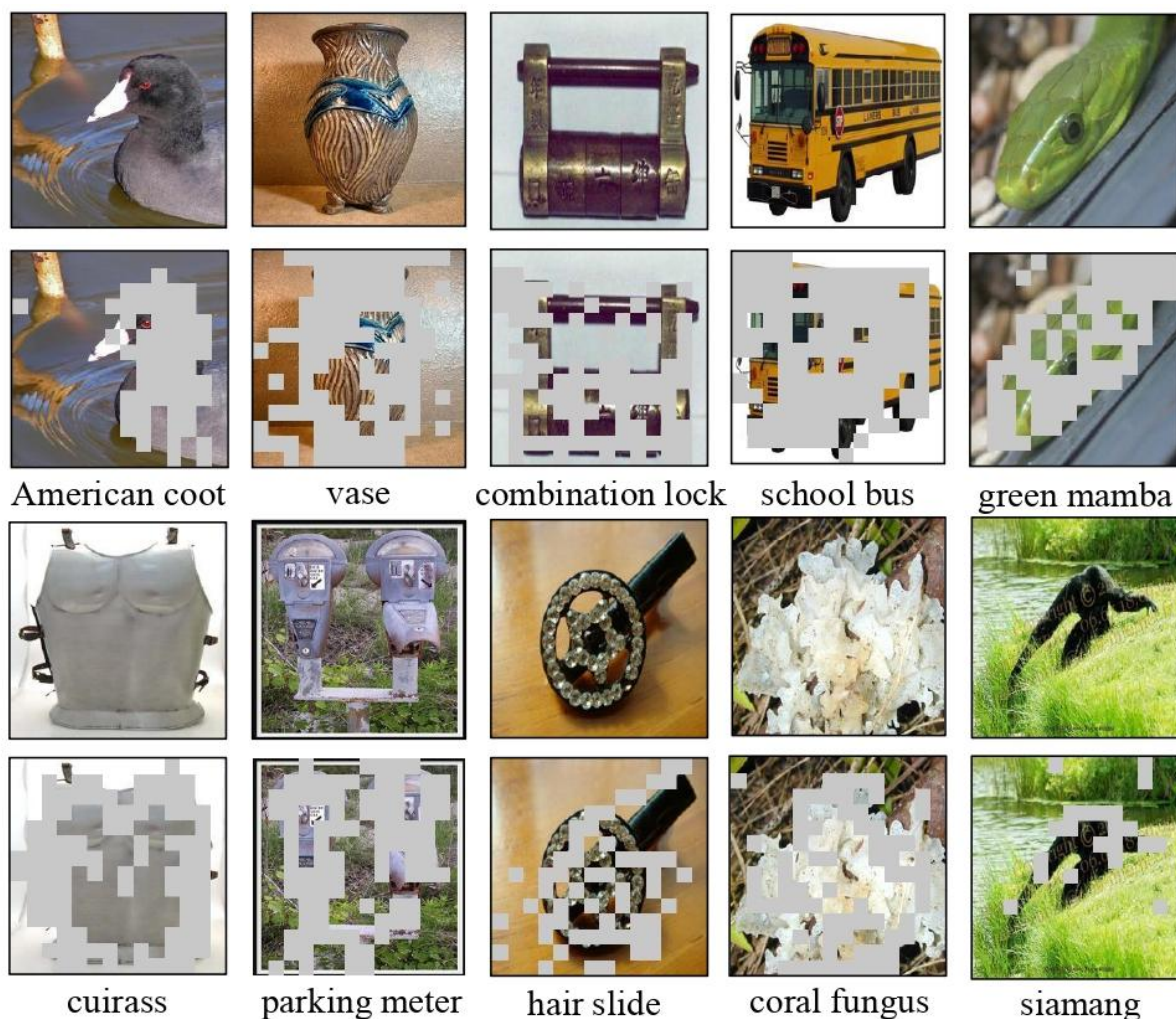


图 2 CLIP 提取的局部特征，灰色区块表示该区块的 Top-1 预测是对应的类别，可以观察到这些区域并没有完美地聚焦于 ID 区域，而是受到了 OOD 背景的严重干扰

然而从 ID 数据中提取的大部分局部上下文并非有效的 OOD 特征。因此，基于这些不可靠 OOD 特征的 OOD 正则化可能会限制 OOD 检测性能的提升。直观地说，当模型对 ID 样本的预测不确定性较高时，从这些数据中提取的 OOD 特征可靠性较低。对这些不可靠的替代 OOD 特征进行正则化会降低 CLIP 的 OOD 检测性能。因此，一个潜在的思路是在模型训练过程中，根据预测不确定性自适应调整从 ID 数据中提取的 OOD 特

征的重要性。基于上述观察，我们提出了一种新的学习框架 SCT，以缓解虚假 OOD 特征引发的问题。总具体来说，我们针对 OOD 检测提示调优的原始学习目标的两个部分，分别引入基于样本不确定性估计的调制因子。在这一新的学习框架下，当使用低置信度数据进行训练时，模型会将注意力转向分类任务，以更好地泛化到下游 ID 数据集；而从高置信度 ID 数据中提取的 OOD 特征则会被赋予更高的权重，以实现更有效的 OOD 正则化。这两个调制因子的重定向作用有助于 CLIP 从非完美的 OOD 特征中学习，最终提升 OOD 检测性能。

同时为了充分利用训练集的视觉信息，本文引入了双模式匹配分数 DPM (Dual-pattern Matching)，通过组合视觉匹配与文本匹配分数（如 MCM 分数 和 GLMCM 分数），在一些应用场景下可以取得更好的检测性能。

本文的创新点如下：

1. 本文探讨了基于 CLIP 的小样本分布外检测问题，通过引入 OOD 正则化和基于不确定度的自适应调制方法，提高了 CLIP 的分布外检测性能。
2. 本文将 DPM 分数与多种小样本提示微调方法结合，在 ImageNet 分布外检测测试基准上的性能有一定的提升，尤其是在 iNaturalist 数据集上。

二、相关工作

近年来，随着机器学习模型在各个领域的广泛应用，其面临的开放环境挑战日益复杂，分布外检测日益成为学界关注的焦点，如新加坡南洋理工大学的研究团队从任务定义和方法分类的角度提出广义分布外检测框架^{[1][2]}，中国科学院自动化研究所的研究团队从问题场景的角度系统综述了分布外检测的最新进展^[3]，本节尝试整合它们的论述思路，对分布外检测的研究现状进行简要的梳理。

（一）单模态分布外检测

单模态分布外检测专注于利用单一图像模态信息判断测试图像是否属于分布外，根据对模型训练的依赖程度，可以划分为事后检测方法 with 训练优化方法两类。

1. 事后检测方法

事后检测方法直接调用训练完成的模型，而无需调整模型参数，即可实现分布外数据的检测，这类方法的核心在于设计判别性良好的评分函数，可分类为基于输出的方法、基于距离的方法、基于梯度的方法、基于特征的方法等。

（1）基于输出的方法：利用模型的输出（如 logit 或 softmax 概率）检测分布外样

本，分布内样本输出一一般较大，而分布外样本输出一一般较小，典型方法有使用最大 softmax 概率的 MSP^[4]、使用最大 logit 的 MLS^[5]、使用能量分数的 EBO^[6]等。

(2) 基于距离的方法：利用特征向量的统计距离度量来检测分布外样本，分布内样本在特征空间中往往形成紧密聚类，样本间距较近，而分布外样本一般与分布内类别的聚类中心距离较远，典型方法有使用马氏距离的 MDS^[7]等。

(3) 基于梯度的方法：利用反向传播的梯度来量化模型的不确定性以检测分布外样本，分布内样本因处于模型的有效学习区域，梯度一般较小且稳定，而分布外样本梯度一般较大且不稳定。典型方法有利用梯度向量范数的 GradNorm^[8]等。

(4) 基于特征的方法：利用特征来检测分布外样本，这是因为分布外样本的特征可能与分布内样本的特征有显著差异，典型方法有对特征添加扰动的 ODIN^[9]，截断中间层高激活值的 ReAct^[10]，利用特征构建虚拟 logit 的 ViM^[11]等。

2. 训练优化方法

训练优化方法旨在训练一个可以对分布内数据准确分类并可以检测出分布外数据的模型，本小节主要讨论基于分类的方法，根据是否使用了分布外数据可以分为无异常值暴露方法和有异常值暴露方法。

(1) 无异常值暴露的方法：典型方法有对 logit 施加恒定向量范数以缓解模型过度自信的 LogitNorm^[12]、将 logit 解耦并平衡各组成部分的 DML^[13]、使用改进的训练目标扩展 ODIN 方法并动态选择扰动幅度的 GODIN^[14]等。

(2) 有异常值暴露的方法：根据异常值的来源可以分为真实异常值暴露方法与伪异常值暴露方法，真实异常值暴露的典型方法有迫使模型将分布外数据优化为均匀分布的 OE^[15]、将分布内数据和分布外数据混合训练的 MixOE^[16]等；而伪异常值暴露的典型方法有从特征空间中类别条件分布的低概率区域采样异常值的 VOS^[17]、基于最近邻的非参数密度估计来选择边界点的 NPOS^[18]等。

(二) 跨模态分布外检测

跨模态分布外检测是指利用图像和文本两种模态信息判断测试图像是否属于分布外。这种方法一般利用视觉语言模型（如 CLIP^[19]等）丰富的图文匹配预训练知识，改善了单模态方法中将标签编码为独热向量造成的语义信息损失问题。

1. 零样本方法

零样本方法不需要在目标任务上进行任何训练，仅依靠预训练模型学习到的知识就能对测试图像进行分布外检测。典型方法有：将视觉文本匹配相似度进行带温度缩放系

数的 softmax 的 MCM^[20]、同时考虑全局匹配分数和局部匹配分数的 GLMCM^[21]、引入分布外类别标签以扩展语义空间的 NegLabel^[22]等。

2. 参数高效微调方法

零样本方法一般不需要对视觉语言模型进行微调,然而其面临与下游数据的语义对齐不足的挑战,而引入一定数量训练样本对视觉语言模型进行参数高效微调的方法有望增强分布外检测方法对下游数据的适应性。

威斯康星-麦迪逊大学的研究团队首次揭示了 CoOp^[23]等参数高效微调方法在分布外检测任务上的巨大潜力^[24]。最近,大量使用参数高效微调方法进行小样本分布外检测的工作大量涌现,如对与分布内无关的局部特征进行熵正则化的 LoCoOp^[25]和 SCT^[26]、利用全局和局部特征来优化全局和局部提示的 GalLop^[27]、尝试为伪分布外样本学习多样性提示的 NegPrompt^[28]和 ID-like^[29]等。

三、方法

(一) 问题重述

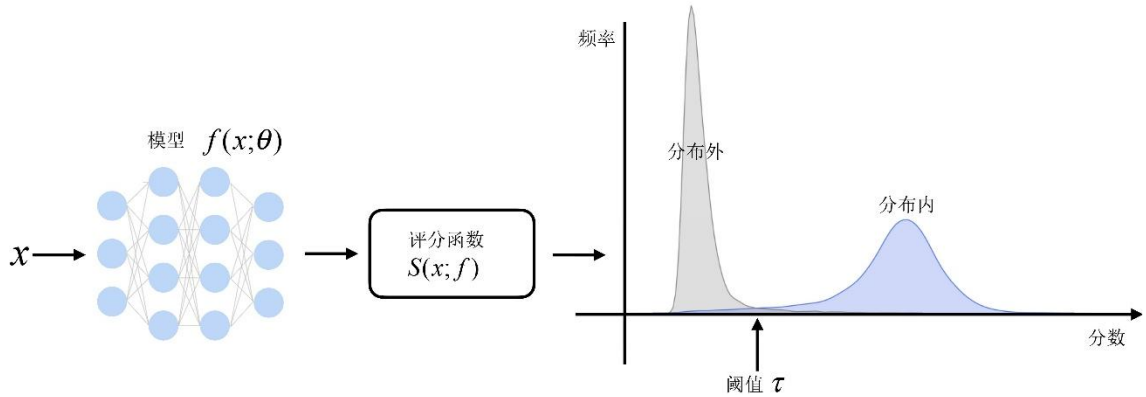


图 3 分布外检测的流程示意图

对于 OOD 检测任务, ID 类别指训练集中使用的类别, OOD 类别是不属于任何 ID 类别的类别。OOD 检测器可视为二分类器, 用于判断图像是 ID 图像还是 OOD 图像。

形式上, 假设我们有一个 ID 数据集 D^{id} , 由 (x^{id}, y^{id}) 对组成, 其中 x^{id} 表示输入的 ID 图像, $y^{id} \in Y^{id} := \{1, \dots, K\}$ 表示 ID 类别标签。令 D^{ood} 表示 OOD 数据集, 由 (x^{ood}, y^{ood}) 对组成, 其中 x^{ood} 表示输入的 OOD 图像, $y^{ood} \in Y^{ood} := \{K + 1, \dots, K + O\}$ 表示 OOD 类别标签。ID 与 OOD 的类别之间没有重叠, 即满足 $Y^{ood} \cap Y^{id} = \emptyset$ 。

本文研究的场景是: 模型仅在训练集 D_{train}^{id} 上进行微调, 且不接触任何真实 OOD 数据。测试集包含 D_{test}^{id} 和 D_{test}^{ood} , 用于评估分布外性能。与现有研究中使用大量 ID 训练样

本或不使用 ID 训练样本不同, 本文的 D_{train}^{id} 是小样本数据集。

(二) CoOp 方法回顾

CoOp 是一种开创性方法, 它利用 CLIP 等视觉-语言预训练知识来实现下游开放世界的视觉识别任务。CLIP 使用手动设计的提示模板, 而 CoOp 则将模板中的部分上下文词设置为可从小样本数据中学习的可学习参数。因此, 分类权重可通过所学提示与视觉特征之间的距离来表示。

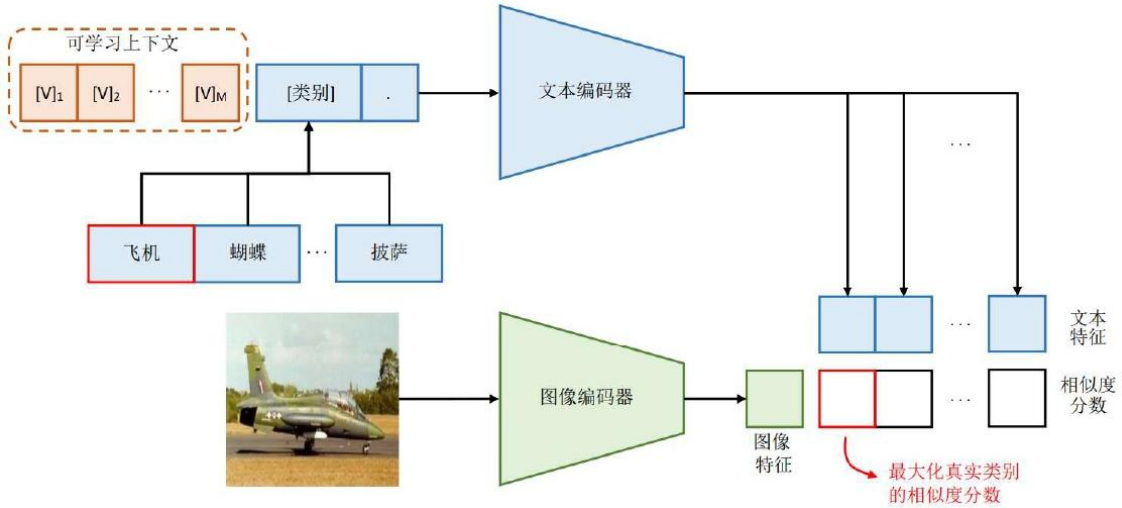


图 4 CoOp 是提示学习的开创性工作

给定一张 ID 图像 x^{id} , 通过 CLIP 的视觉编码器 Φ^V 获得全局视觉特征 $F^{id} = \Phi^V(x^{id})$ 。然后, 文本提示可表示为 $t_M = \{\omega_1, \omega_2, \dots, \omega_N, c_m\}$, 其中 c_m 表示 ID 类别名称的词嵌入, $\omega = \{\omega_n | n=1\}^N$ 为可学习向量 (每个向量与原始词嵌入维度相同, N 表示上下文词的长度)。以提示 t_m 为输入, 文本编码器 Φ^T 输出文本特征 $F_m^T = \Phi^T(t_m)$ 。最终预测概率通过匹配分数计算如下:

$$p(y = m | x^{id}) = \frac{\exp(\text{sim}(F^{id}, F_m^T)/\tau)}{\sum_{j=1}^M \exp(\text{sim}(F^{id}, F_j^T)/\tau)}$$

其中 $\text{sim}(\cdot, \cdot)$ 表示余弦相似度, τ 为 Softmax 温度参数。

最终损失 \mathcal{L}_{coop} 为公式 (1) 与真实标签 y^{in} 的交叉熵损失。

$$\mathcal{L}_{coop} = \mathcal{L}_{CE}(p(y|x), y)$$

(三) OOD 正则化

在本节中, 我们介绍用于小样本 OOD 检测的 LoCoOp 方法。图 5 展示了本文提出的方法的总体框架。该方法包含两个组成部分: 第一部分是从 CLIP 局部特征中提取与 ID 无关的区域, 第二部分是利用提取的区域进行 OOD 正则化训练。另外为了平衡 ID

分类损失与 OOD 正则化损失，本文还介绍了对损失项的调制方法。

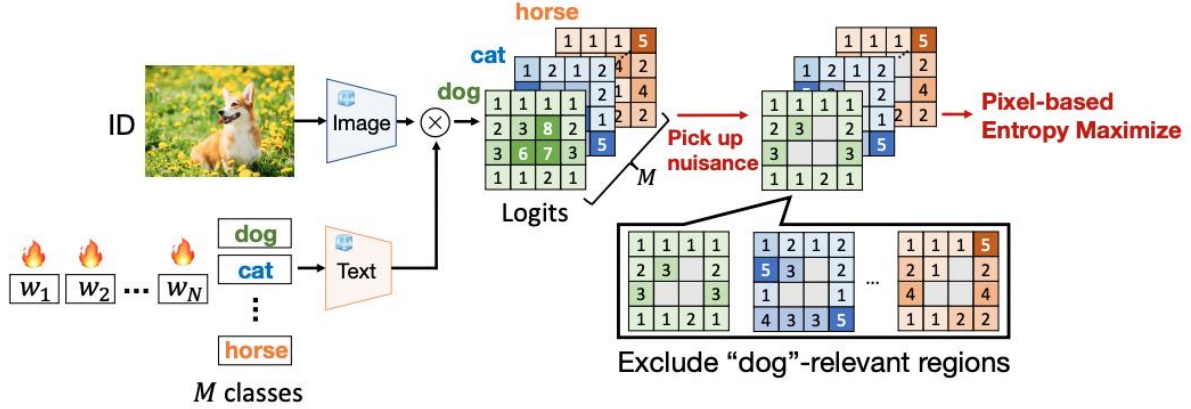


图 5 本文的 LoCoOp 及 SCT 方法的总览图

下面将简要概述如何从 CLIP 中获取局部特征。然后我们将描述两个关键组件：从 ID 图像中提取与目标无关区域的方法，以及训练过程中使用的 OOD 正则化损失。

(1) 从 CLIP 中获取局部特征

为了获取 CLIP 的局部特征，本文按如下方式将每个区域 i 的视觉特征图特征 F_i^{id} 投影到文本空间：

$$F_i^{id} = \text{Proj}_{v \rightarrow t}(v(F_i^{id})),$$

其中 v 表示值投影， $\text{Proj}_{v \rightarrow t}$ 表示 CLIP 中从视觉空间到文本空间的投影，该局部特征具有丰富的局部视觉-文本对齐信息。

(2) ID 无关区域的提取

我们需要从所有区域索引集合 $I = \{0, 1, 2, \dots, H \times W - 1\}$ 中选择 ID 无关区域的索引，其中 H 和 W 分别表示特征图的高度和宽度。我们以一种简单直观的方式处理这一问题：与分割任务类似，在训练过程中，通过计算每个区域 i 的图像特征 F_i^{id} 与 ID 类别的文本特征之间的相似度，可得到每个区域 i 的分类预测概率，公式如下：

$$p_i(y = m | x^{id}) = \frac{\exp(\text{sim}(F_i^{id}, F_m^T)/\tau)}{\sum_{j=1}^M \exp(\text{sim}(F_i^{id}, F_j^T)/\tau)}$$

当 x^{id} 的区域 i 对应 ID 目标的一部分时，真实类别 y^{id} 应位于前 K 个预测类别中。相反，若区域 i 与 ID 无关（如背景），由于该区域与真实标签 y^{id} 缺乏语义关联， y^{id} 不会出现在前 K 个预测类别中。基于此观察，我们将真实类别未包含在前 K 个预测类别中的区域识别为 ID 无关区域 J ，公式如下：

$$J = \{i \in I : \text{rank}(p_i(y = y^{id} | x^{id})) > K\}$$

其中 $\text{rank}(p_i(y = y^{id}|x^{id}))$ 表示真实类别 y^{id} 在所有 ID 类别中的预测排名。需要注意的是, 图像 x^{id} 的区域集合 J 会在训练过程中随 $p_i(y = k|x^{id})$ 的更新而动态调整。

这种方法可能看起来是一种依赖参数的方法, 因为它依赖于超参数 K 。然而我们认为最优的 K 并不难搜索, 这是因为在现实世界中应用 OOD 检测时, 我们假设用户了解什么是 ID 数据。例如, 当 ID 是 ImageNet-1K 时, 用户事先知道细粒度类别的数量或类别之间的语义关系, 利用细粒度类别数量等先验知识, 有助于确定 K 的值。

(3) OOD 正则化

利用 ID 无关区域的索引集合 J , 我们进行 OOD 正则化, 以从文本特征中剔除不必要的信息。每个属于 J 的 OOD 区域 j 的特征 F_j^{id} 应与任何 ID 文本嵌入不相似。因此, 我们采用熵最大化进行正则化——该方法在训练中常用于检测未知样本。熵最大化使 $p_j(y|x^{id})$ 的熵值更大, 从而让模型能够 OOD 图像特征 F_j^{id} 与任何 ID 文本嵌入不相似, 因此该正则化的损失函数如下:

$$\mathcal{L}_{ood} = H(p_j) = -\frac{1}{N} \sum_{j=1}^N p_j \log p_j$$

其中, $H(\cdot)$ 为熵函数, p_j 表示属于集合 J 的区域 j 的预测概率, N 表示集合 J 的大小。

LoCoOp 的损失函数可以表示如下, 其中 λ 为超参数:

$$\mathcal{L} = \mathcal{L}_{coop} + \lambda \cdot \mathcal{L}_{ood}$$

(四) 基于不确定度的自适应调制

如前所述, 基于 LoCoOp 的 OOD 检测范式依赖于从 ID 数据中提取的背景局部特征进行 OOD 正则化。然而 ID 前景与 OOD 背景的分解并不完美, 模型需要从这些不准确的 OOD 特征中有效学习以实现更好的 OOD 检测。缓解这一问题的一个思路是根据不确定性估计, 自适应地调整从不同 ID 样本生成的 OOD 正则化的重要性, 从而减轻无效 OOD 特征的错误引导。在这种学习范式下, 模型可以由更有效的 OOD 特征进行正则化, 同时避免过度自信, 进而提升 OOD 检测性能。基于这一思路, 本文考虑在提示调优的框架下, 将学习目标重新表述如下:

$$\mathcal{L}_{SCT} = \mathcal{L}_{coop} \cdot \varphi(p(y|x)) + \lambda \cdot \mathcal{L}_{ood} \cdot \psi(p(y|x))$$

其中, $\varphi: \mathbb{R}^M \rightarrow \mathbb{R}$ 和 $\psi: \mathbb{R}^M \rightarrow \mathbb{R}$ 表示新引入的调制函数, 它们根据不确定性估计, 为 LoCoOp 原始损失函数的两个分量计算自适应因子。在这个损失函数中, 左边部分用于 ID 分类任务, 右边部分用于 OOD 正则化。具体而言, φ 关于 $p(y|x)$ 应单调递减, ψ

关于 $p(y|x)$ 应单调递增，这样在训练过程中，调制因子就能在两个任务之间调整提示学习的重点。当模型对真实标签输出低置信度的预测时，ID 分类任务的重要性会被凸显出来，以便更好地泛化到下游任务，同时减少从 ID 数据中提取的无效 OOD 特征的正则化影响。当模型能够准确且自信地对 ID 样本进行分类时，其注意力会转向 OOD 正则化，以强化有用的与 ID 无关特征的积极作用，从而实现更好的 OOD 检测。与此同时，分类任务的损失贡献会降低，以避免模型过度拟合下游数据集，这有利于模型的校准，并进一步提高提取的 OOD 特征的有效性。

在满足上述简单要求的众多函数中，我们选择设计简单的线性调制函数，因为其设计简单。具体而言，我们将 SCT 的损失函数公式定义如下：

$$\mathcal{L}_{SCT} = \mathcal{L}_{coop} \cdot (1 - p(y|x)) + \lambda \cdot \mathcal{L}_{ood} \cdot p(y|x)$$

(五) 测试时 OOD 检测分数

在测试阶段，我们采用 MCM 分数、GL-MCM 分数以及 DPM 分数。通过这些函数，我们将测试样本分类为 ID 图像 x^{id} 和 OOD 图像 x^{ood} 。具体细节如下：

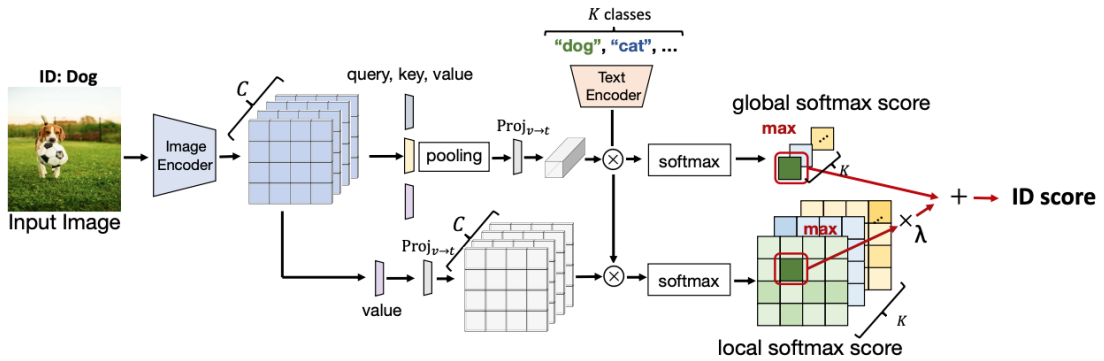


图 6 MCM 分数与 GLMCM 分数计算示意图

(1) MCM 分数

利用全局图像特征和文本的 softmax 分数。其基本原理是：对于 ID 数据，它会以高置信度匹配到某个文本原型，而对于 OOD 数据则相反，其形式化表达为：

$$S_{MCM} = \max_m \frac{\exp(\text{sim}(F, F_m^T)/\tau)}{\sum_{j=1}^M \exp(\text{sim}(F, F_j^T)/\tau)}$$

其中 τ 是温度缩放系数。

(2) GL-MCM 分数

利用全局和局部图像特征以及文本特征的 softmax 分数，其形式化表达为：

$$S_{GLMCM} = S_{MCM} + \lambda \cdot \max_{m,i} \frac{\exp(\text{sim}(F_i, F_m^T)/\tau)}{\sum_{j=1}^M \exp(\text{sim}(F_i, F_j^T)/\tau)}$$

当 ID 图像中出现 OOD 对象时，全局图像特征的匹配分数可能会被错误地降低，利用局部特征的匹配分数可以补偿较低的全局匹配分数，并产生正确的 ID 置信度。

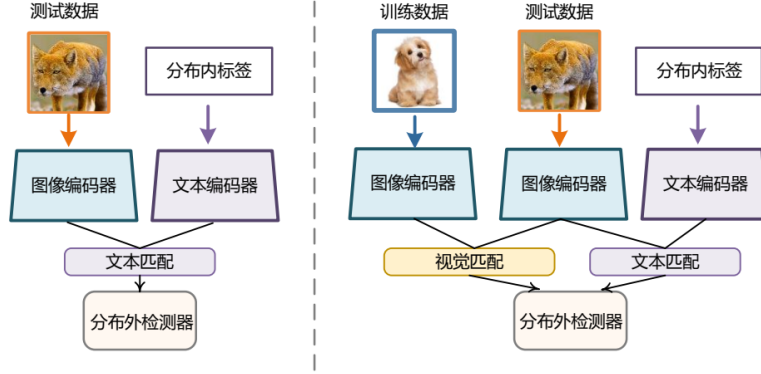


图 7 DPM 分数计算示意图

(3) DPM 分数

首先定义文本模式与视觉模式，对于一个具有 M 个类别的任务，类别 m 的文本模式是文本特征 $P_m^T = F_m^T = \Phi^T(t_m) \in \mathbb{R}^{1 \times D}$ ，将类别 m 的第 j 个样本记为 x_j^k ，其图像-文本匹配分数 s_j^m 为：

$$s_j^m = \text{Concat}[p^T(y_1|x_j^m), p^T(y_2|x_j^m), \dots, p^T(y_K|x_j^m)] \in \mathbb{R}^{1 \times M}$$

将类别 k 的图像-文本匹配分数的均值作为类别 m 的视觉模式 P_k^V ，即

$$P_m^V = \frac{1}{N_m} \sum_{j=1}^{N_m} s_j^m \in \mathbb{R}^{1 \times M}$$

其中 N_m 是类别 m 的样本数。

现在我们有视觉模式 P^V 和文本模式 P^T ，如下所示：

$$P^V = \text{Concat}([P_1^V, P_2^V, \dots, P_M^V]) \in \mathbb{R}^{M \times M}$$

$$P^T = \text{Concat}([P_1^T, P_2^T, \dots, P_M^T]) \in \mathbb{R}^{M \times D}$$

对于一个待识别的图像 x ，我们可以计算其与文本模式的匹配分数：

$$s_m^T = p^T(y_m|x) \in \mathbb{R}$$

$$s^T = \text{Concat}([s_1^T, s_2^T, \dots, s_M^T]) \in \mathbb{R}^{1 \times M}$$

考虑到视觉语言模型中存在固有的模态差距，因此使用 KL 散度来度量测试特征与视觉模式之间的相似度，其计算公式为 $KL(R \parallel Q) = \sum_i R_i \log \frac{R_i}{Q_i}$ ，我们可以计算其文本原型引导的聚合特征与视觉模式的匹配分数：

$$s_m^V = KL(s^T \parallel P_m^V) \in \mathbb{R}$$

$$s^V = \text{Concat}([s_1^V, s_2^V, \dots, s_M^V]) \in \mathbb{R}^{1 \times M}$$

每个模式匹配都强调并贡献来自不同模态的信息，因此利用这两个匹配分数来融合

来自两种模态的信息。形式上, 该分布外检测评分函数 DPM 可以表示为:

$$DPM(x) = \max(s^T) - \beta \min(s^V) \in \mathbb{R}, \quad OOD = \begin{cases} 1, & \text{if } DPM(x) > \lambda \\ 0, & \text{if } DPM(x) < \lambda \end{cases}$$

其中 β 是视觉模态亲和因子, 用于灵活地平衡视觉和文本信息, 在实际应用中一般会选择 λ 使得高比例的分布内数据 (例如, 95%) 高于阈值。

在计算过程中, 两种匹配分数可能会存在取值范围的差异, 因此将视觉匹配分数缩放到文本匹配分数的取值范围, 计算方式如下:

$$s^V = \frac{s^V - \min(s^V)}{\max(s^V) - \min(s^V)} \cdot (\max(s^T) - \min(s^T)) + \min(s^T)$$

四、数据

遵循 MOS^[30]中的任务设置, 我们采用如下数据集。



图 8 MOS 中提出的大规模图像分类分布外检测测试基准, 以 ImageNet-1K 作为分布内数据集, iNaturalist、SUN、Places、Texture 作为分布外数据集

(一) 分布内数据集

ImageNet-1K 数据集^[31]: 该数据集为 ILSVRC 2012 竞赛数据集, 为完整的 ImageNet 数据集一个较为知名的子数据集, 包含 1000 个类别, 每个类别有 1000 余个样本, 总计 120 万训练样本和 5 万测试样本, 其组织结构遵循严格的层次分类体系, 确保每个类别的图像具有高度的代表性和多样性。该数据集图像质量高, 分辨率多样, 涵盖了从自然场景到人工制品的广泛领域, 以其庞大的规模和丰富的类别多样性著称, 为图像分类、域适应和半监督学习等任务提供了坚实的基础。

(二) 分布外数据集

(1) iNaturalist 数据集^[32]: 该细粒度数据集包含 859,000 张图像, 涵盖 5,000 多种动植物, 所有图像都被调整大小, 最大尺寸为 800。研究人员手动挑选了 110 个在 ImageNet-1K 中不存在的植物类别, 然后为这 110 个类别中随机采样 10,000 张图像。

(2) SUN 数据集^[33]: 该数据集包含 397 个类别、130,519 张 200×200 以上分辨率

的场景图片。研究人员手动挑选了 50 个在 ImageNet-1K 中不存在的自然场景，如森林、冰山等，然后为这 50 个类别随机采样 10,000 张图像。

(3) Places 数据集^[34]：该数据集是另一个与 SUN 概念覆盖相似的场景数据集，这个数据集中的所有图像都被调整大小，最小尺寸为 512。研究人员手动挑选 50 个 ImageNet-1K 中不存在的类别，然后为这 50 个类别随机采样 10,000 张图像。

(4) Texture 数据集^[35]：该数据集由 5,640 张纹理图案图像组成，尺寸范围在 300×300 到 640×640 之间，由于与 ImageNet-1K 没有类别重叠，因此使用整个数据集。

对这些数据集的特征进行 t-SNE 降维可视化，如图 9 所示，可以看出 iNaturalist 与 ID 的重叠最少，其次是 SUN 和 Places，而 Texture 与 ID 的重叠最多。

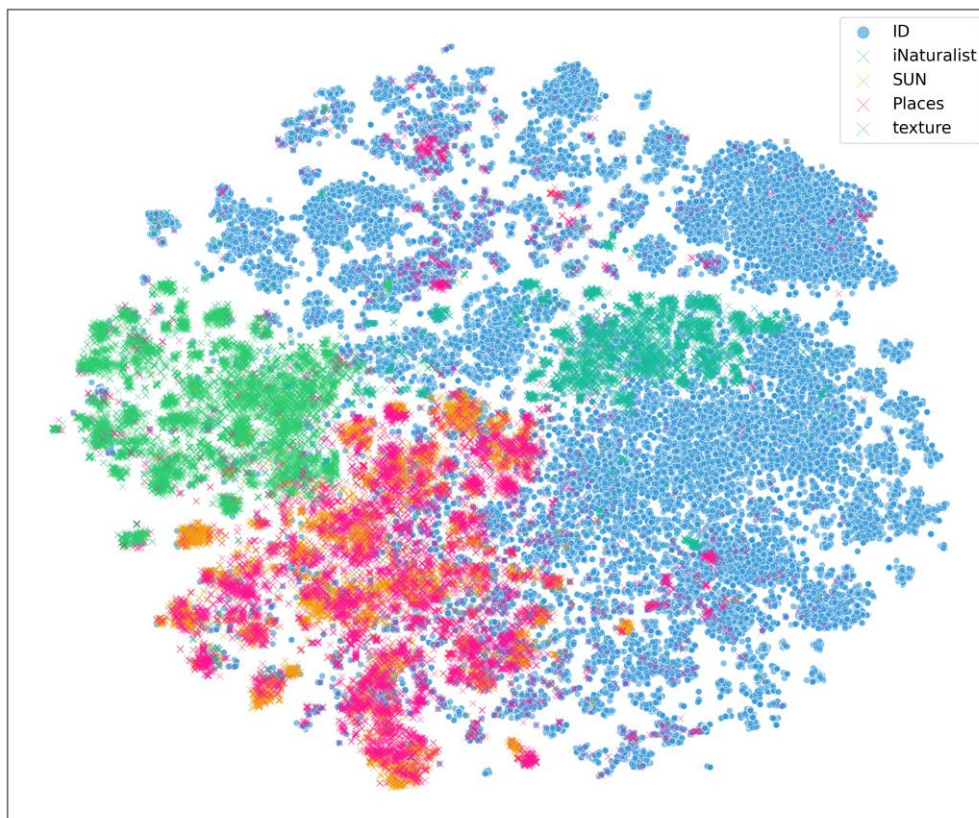


图 9 ID 数据集和 OOD 数据集 t-SNE 特征降维可视化

五、实验结果

(一) 实验设置

本文所有 CLIP 小样本微调方法均在单张 Nvidia 4090 上运行，采用小样本学习的设置，每个类别采样 1, 2, 4, 8, 16，其中计算开销最大的 16-shot 训练所需内存约 20GB，推理所需内存约 4GB，训练参数如下表，其他更具体的运行配置详见 github 的 README。

表 1 实验参数设置

参数	取值
学习率	0.002
可学习提示的长度	16
优化器	SGD
学习率调度器	余弦退火
预热训练轮次	1
预热学习率	0.0001
数据增强	RandomFlip、RandomResizedCrop
批次大小	32

(二) 评估指标

本文使用以下指标来衡量分布外检测的性能：

- (1) ID 样本的真阳性率为 95% 时 OOD 样本的假阳性率($FPR@95$, 可简写为 FPR)。
- (2) 接收者工作特征曲线下面积 ($AUROC$, 可简写为 AUR)。

(三) 对比方法

本文选取了近年来多种具有代表性的方法进行对比实验，这些对比方法基本囊括了当前分布外检测领域的主流技术方向：

(1) 单模态方法

- ①EBO^[6](NeurIPS 2020)，性能报道来自 NPOS 论文。
- ②ODIN^[9](NeurIPS 2021)，性能报道来自 NPOS 论文。
- ③ViM^[11](CVPR 2023)，性能报道来自 NPOS 论文。
- ④VOS^[17] (ICLR 2022)，性能报道来自 NPOS 论文。
- ⑤NPOS^[18] (ICLR 2023)，性能报道来自 NPOS 论文。

(2) CLIP 零样本方法

- ①MCM^[20] (NeurIPS 2022)，性能报道来自 LoCoOp 论文。
- ②GL-MCM^[21] (IJCV 2025)，性能报道来自 LoCoOp 论文。
- ③CLIPN^[36] (ICCV 2023)，性能报道来自 CLIPN 论文。

(3) CLIP 小样本微调方法

- ①CoOp^{[23][24]} (IJCV 2021)，正文报道复现结果，附录报道 LoCoOp 论文的结果。
- ②LoCoOp^[25] (NeurIPS 2023)，正文报道复现结果，附录报道 LoCoOp 论文的结果。
- ③SCT^[26] (NeurIPS 2024)，正文报道复现结果，附录报道 SCT 论文的结果。

上述 CLIP 小样本微调方法在论文中报道的性能见附表 6，各种小样本设置下三个随机数种子的具体实验结果见附表 1、附表 2、附表 3、附表 4、附表 5。

(四) 对比实验结果

表 2 各种分布外检测方法在 ImageNet-1K OOD 测试基准上的性能对比

对比方法	iNaturalist		SUN		Places		Texture		Average	
	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓
单模态方法										
EBO	94.68	29.75	87.33	53.18	85.60	56.40	88.00	51.35	88.90	47.67
ODIN	94.65	30.22	87.17	54.04	85.54	55.06	87.85	51.67	88.80	47.75
ViM	93.16	32.19	87.19	54.01	83.75	60.67	87.18	53.94	87.82	50.20
VOS+	94.62	28.99	92.57	36.88	91.23	38.39	86.33	61.02	91.19	41.32
NPOS	96.19	16.58	90.44	43.77	89.44	45.27	88.80	46.12	91.22	37.93
CLIP 零样本方法										
MCM	94.61	30.94	92.56	37.67	89.76	44.76	86.10	57.91	90.76	42.82
GL-MCM	96.71	15.18	93.09	30.42	89.90	38.85	83.63	57.93	90.83	35.47
CLIPN	95.27	23.94	93.93	26.17	90.93	40.83	92.28	33.45	93.10	31.10
CLIP 小样本微调方法 (1-shot)										
CoOp	91.96	41.56	92.25	37.61	89.70	44.15	87.60	51.19	90.38	43.63
LoCoOp	94.79	25.72	94.70	24.27	91.73	33.24	87.56	51.48	92.20	33.68
SCT	94.75	25.20	94.00	26.93	91.33	34.44	86.62	52.01	91.68	34.65
DPCoOp	94.33	29.18	92.25	37.61	89.70	44.15	87.65	51.07	90.98	40.50
DPLoCoOp	96.30	17.71	94.70	24.27	91.73	33.24	87.74	51.22	92.62	31.61
DPST	96.50	16.49	94.00	26.93	91.33	34.44	86.83	51.67	92.17	32.38
CLIP 小样本微调方法 (2-shot)										
CoOp	92.78	35.25	92.23	35.07	89.74	42.68	89.53	44.89	91.07	39.47
LoCoOp	96.40	16.64	95.11	22.90	92.18	32.05	89.97	43.75	93.42	28.84
SCT	95.65	20.72	94.57	24.40	91.56	33.19	87.74	48.75	92.38	31.77
DPCoOp	95.44	21.92	92.23	35.07	89.74	42.68	89.57	44.69	91.75	36.09
DPLoCoOp	97.32	12.36	95.11	22.90	92.18	32.05	90.13	43.40	93.69	27.68
DPST	97.31	12.12	94.57	24.40	91.56	33.19	88.05	47.98	92.87	29.42
CLIP 小样本微调方法 (4-shot)										
CoOp	93.52	31.33	92.68	34.65	89.91	42.33	89.37	44.61	91.37	38.23
LoCoOp	96.20	17.36	95.33	22.27	92.11	32.10	89.37	43.89	93.25	28.91
SCT	96.33	16.71	95.33	20.86	92.21	30.41	88.83	44.12	93.18	28.03
DPCoOp	95.38	22.38	92.68	34.65	89.91	42.33	89.50	44.50	91.87	35.97
DPLoCoOp	97.13	13.10	95.33	22.27	92.11	32.10	89.61	43.59	93.55	27.77
DPST	97.56	10.51	95.33	20.86	92.21	30.41	89.02	44.02	93.53	26.45
CLIP 小样本微调方法 (8-shot)										
CoOp	93.80	30.96	93.39	30.81	90.64	39.35	89.61	45.03	91.86	36.54
LoCoOp	95.73	20.49	95.41	21.60	92.42	31.21	89.78	45.14	93.34	29.61
SCT	96.16	17.92	94.96	22.33	92.27	31.20	89.22	43.95	93.15	28.85
DPCoOp	95.75	20.90	93.39	30.81	90.64	39.35	89.66	45.04	92.36	34.03
DpLoCoOp	97.04	13.52	95.41	21.60	92.42	31.21	89.98	45.00	93.71	27.83
DPST	97.45	11.44	94.96	22.33	92.27	31.20	89.45	43.50	93.53	27.12
CLIP 小样本微调方法 (16-shot)										
CoOp	94.23	28.78	92.67	34.77	89.93	42.74	89.74	44.33	91.64	37.66
LoCoOp	96.44	16.81	94.90	23.27	92.16	32.24	90.42	41.47	93.48	28.45
SCT	96.32	16.27	95.22	21.52	92.43	29.83	89.37	42.87	93.34	27.62
DPCoOp	96.33	17.08	92.67	34.77	89.93	42.74	89.78	44.13	92.18	34.68
DPLoCoOp	97.56	10.99	94.90	23.27	92.16	32.24	90.55	41.27	93.79	26.94
DPST	97.72	9.80	95.22	21.52	92.43	29.83	89.45	42.83	93.71	26.00

表 2 展示了单模态方法、CLIP 零样本方法、CLIP 小样本微调方法（1-shot、2-shot、4-shot、8-shot、16-shot）在 ImageNet-1K OOD 测试基准上的性能对比，可以看出基于 CLIP 的方法普遍优于单模态的方法，而 CLIP 小样本微调方法普遍优于 CLIP 零样本方法，且随着样本数的增加性能有普遍提升。LoCoOp 和 SCT 在引入了 OOD 正则化之后相较于 CoOp 方法有比较显著的提升，所有小样本微调方法在使用了 DPM 分数之后性能有了进一步的提升，尤其是在 iNaturalist 这个 OOD 数据集上。但是 SCT 相比 LoCoOp 性能没有明显的提升，这表明基于不确定性的调制函数的有效性值得进一步商榷。值得注意的是，使用 DPM 分数后，三种方法在 SUN 和 Places 这两个场景数据集上并没有提升，这表明了模型仍然很难区分纯背景的 OOD 样本，值得进一步讨论与改进。

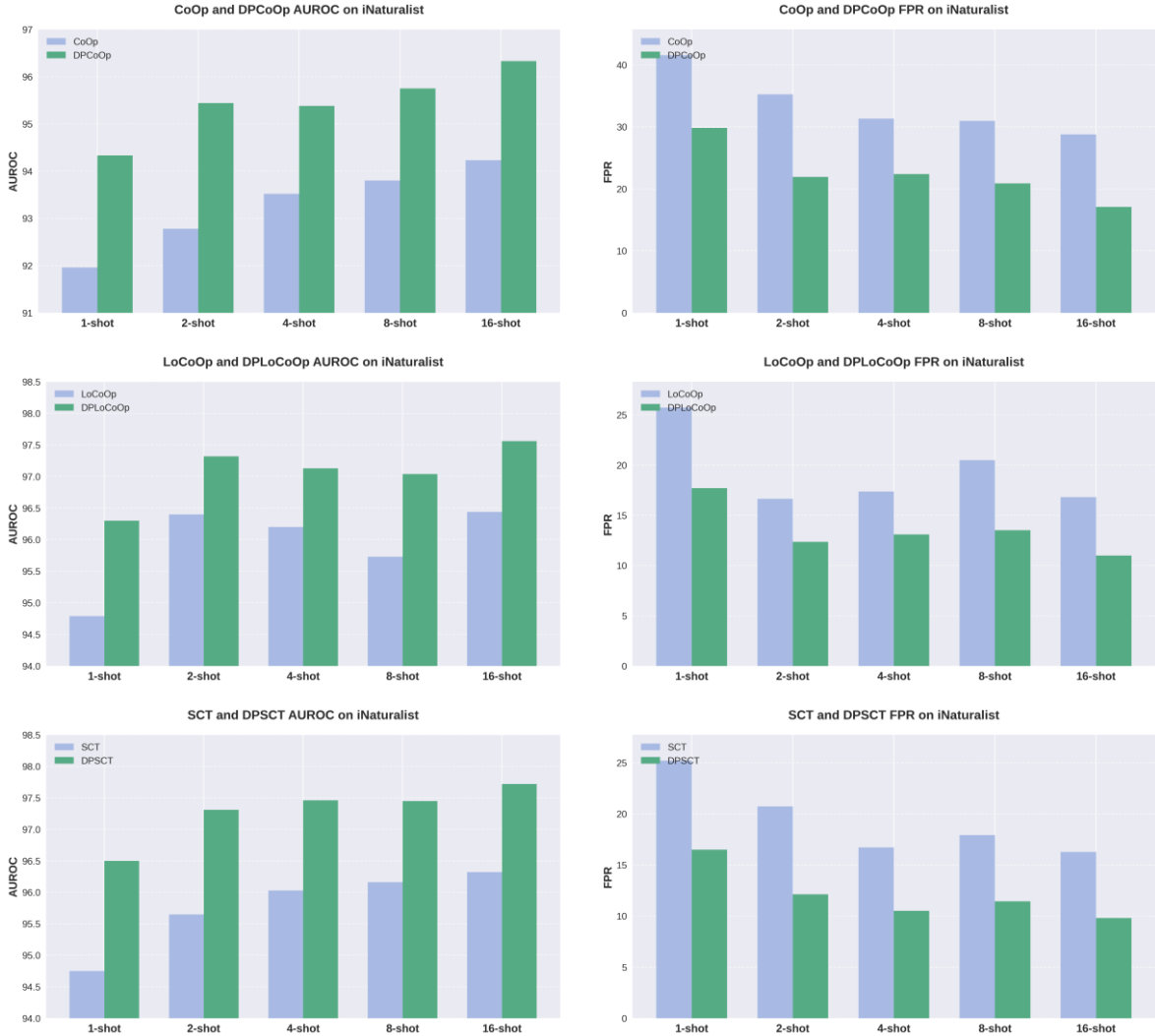


图 10 三种方法使用 DPM 分数后在 iNaturalist 数据集上的 AUR 和 FPR 性能提升对比

图 10 展示了三种方法使用 DPM 分数后在 iNaturalist 数据集上的 AUR 和 FPR 性能提升对比，可以看出使用了 DPM 分数之后 AUROC 显著提高，FPR@95 显著降低，这

表明 ID 与 OOD 的可分离性增强。图 11 展示的 CoOp 方法和 DPCoOp 方法在 1-shot 下两个不同的随机数种子下在 iNaturalist 数据集上的检测分数分布图进一步证明了这一点。

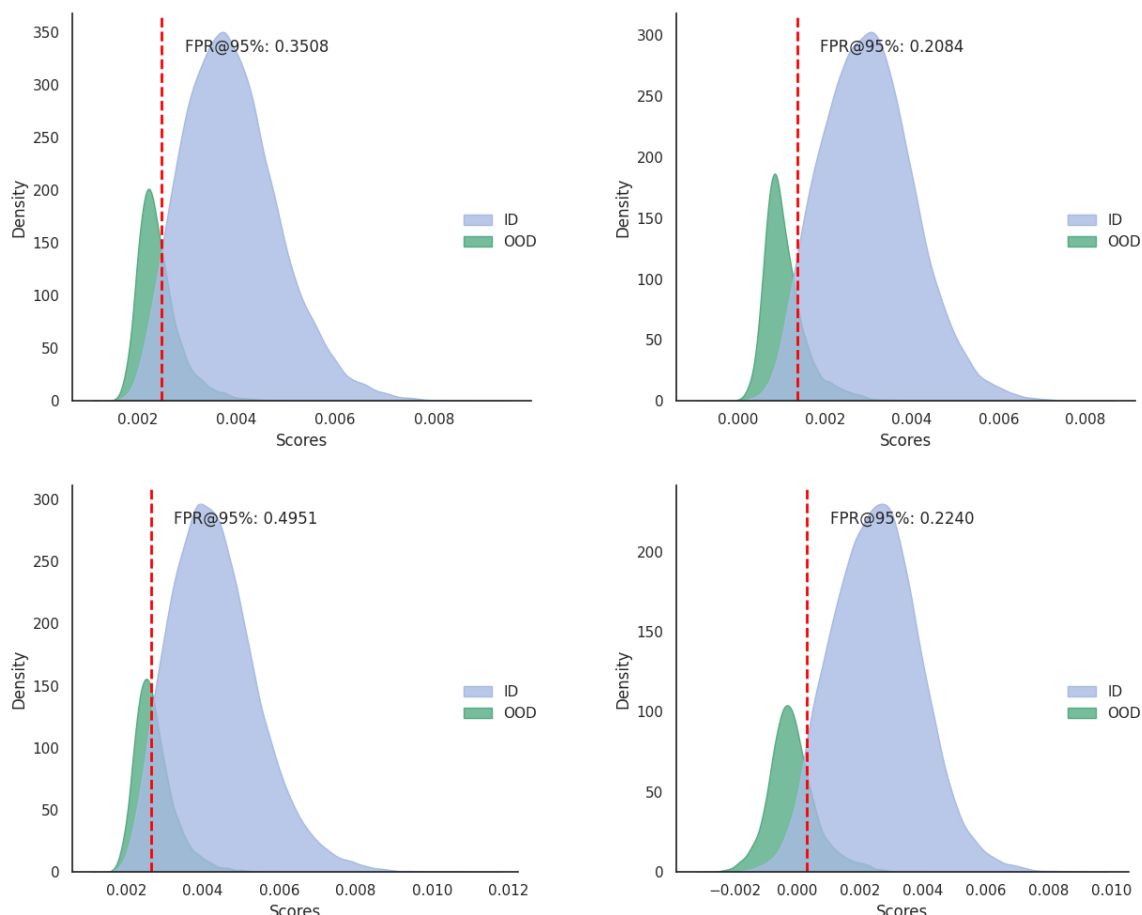


图 11 CoOp 方法和 DPCoOp 方法在 1-shot 下两个不同随机数种子下在 iNaturalist 数据集上的检测分数分布图，左：CoOp 方法，右：DPCoOp 方法

表 3 展示了不同方法在 ImageNet-1K 测试集上的分类准确率们可以看出 CLIP 方法总体不如单模态方法，而在 CLIP 方法中，随着样本数量的增加，三种提示学习方法的准确率稳步提升，同时值得注意的是，SCT 方法在所有情况下都取得了最好的结果，这表明本文设计的调制函数确实平衡了两个任务的权重，这种改进有效提高了分类准确率。

表 3 各类方法在 ImageNet-1K 测试集上的分类准确率

对比方法	1-shot	2-shot	4-shot	8-shot	16-shot
单模态方法					
Posthoc			79.64		
VOS+			79.58		
NPOS			79.42		
CLIP 方法					
Posthoc			66.70		
CoOp	69.08	69.31	70.24	70.88	71.20
LoCoOp	69.14	69.60	70.03	70.73	71.05
SCT	69.25	69.78	70.37	70.97	71.48

(五)消融实验结果

关于 OOD 损失与 DPM 分数的消融实验，在表 2 已经有较好的呈现。

表 4 展示了 SCT 方法中线性调制函数在 1-shot 和 16-shot 这两个场景下的消融实验，其中在 1-shot 场景下只保留分类损失前的调制函数可以取得最佳效果，16-shot 场景下同时在分类损失和 OOD 损失前使用调制函数可以取得最佳性能。

表 4 SCT 线性调制函数的消融实验结果

φ	ψ	iNaturalist		SUN		Places		Texture		Average	
		AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓
小样本微调（1-shot）											
×	×	94.05	28.81	94.51	25.76	91.59	33.68	86.85	51.53	91.75	34.94
√	×	95.14	21.94	95.08	22.46	92.04	30.54	87.14	49.61	92.35	31.14
×	√	95.72	19.26	94.44	23.70	91.04	33.56	85.77	51.10	91.74	31.90
√	√	95.70	19.16	94.58	23.52	91.23	32.81	86.66	48.87	92.04	31.09
小样本微调（16-shot）											
×	×	96.30	17.58	95.20	22.82	92.03	32.21	88.86	45.27	93.10	29.47
√	×	94.50	18.14	94.04	22.42	91.18	31.90	90.09	44.72	92.66	29.30
×	√	96.92	15.08	95.16	21.42	92.07	30.60	86.35	48.64	92.62	28.94
√	√	95.86	13.94	95.33	20.55	92.24	29.86	89.06	41.51	93.37	26.47

表 5 SCT 不同的调制函数

调制函数	φ	ψ
线性函数 (SCT-L)	$1 - p(y \mathbf{x})$	$p(y \mathbf{x})$
幂函数 (SCT-P)	$(1 - p(y \mathbf{x}))^\alpha$	$p(y \mathbf{x})^\alpha$
对数函数 (SCT-Log)	$1 - \frac{\log(p(y \mathbf{x}) + 1)}{\log 2}$	$\frac{\log(p(y \mathbf{x}) + 1)}{\log 2}$
三角函数 (SCT-Tri)	$\cos(\frac{\pi}{2}p(y \mathbf{x}))$	$\sin(\frac{\pi}{2}p(y \mathbf{x}))$

为了探索本文提出的学习框架的通用性，我们考虑了调制函数的其他实例。具体来说，我们考虑了幂函数、对数函数和三角函数，其公式见表 5。在幂函数实验中，我们设置 1-shot 场景下 $\lambda = 0.25$ 、 $\alpha = 0.5$ ，16-shot 场景下 $\lambda = 0.25$ 、 $\alpha = 4$ 。如表 6 所示的实验结果表明，本文所采用的线性调制函数在两种小样本场景下都能取得最优的 FPR 性能，但是对于这些调制函数对性能影响更理论性的分析还有待完善。

表 6 SCT 使用不同调制函数的消融实验结果

对比 方法	iNaturalist		SUN		Places		Texture		Average	
	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓
小样本微调 (1-shot)										
SCT-L	95.70	19.16	94.58	23.52	91.23	32.81	86.66	48.87	92.04	31.09
SCT-Pow	95.84	18.91	94.61	25.06	91.68	33.36	86.68	49.13	92.20	31.62
SCT-Log	95.98	18.32	94.62	26.00	91.39	33.42	85.97	51.17	91.99	32.23
SCT-Tri	95.84	18.46	94.98	24.32	91.69	33.23	86.15	54.10	92.17	32.53
小样本微调 (16-shot)										
SCT-L	95.86	13.94	95.33	20.55	92.24	29.86	89.06	41.51	93.37	26.47
SCT-Pow	96.70	14.07	94.75	20.74	91.96	30.11	87.90	43.49	92.83	27.10
SCT-Log	97.01	13.11	95.50	20.56	92.66	29.03	87.61	45.55	93.20	27.06
SCT-Tri	96.81	14.88	95.44	20.30	92.50	29.33	87.90	44.84	93.16	27.34

六、总结

本报告聚焦基于 CLIP 的小样本分布外 (OOD) 检测问题, 提出 LoCoOp 和 SCT 方法提升检测性能。LoCoOp 通过提取 ID 数据中的 OOD 背景特征进行正则化, 抑制干扰以增强 ID 与 OOD 区分度。针对其不足, SCT 引入调制因子, 根据预测不确定性自适应校准 OOD 正则化影响, 同时引入双模式匹配分数 DPM 融合多模态信息。实验表明, LoCoOp 和 SCT 在 ImageNet 基准上性能优于对比方法, 结合 DPM 分数性能有进一步提升, 尤其在 iNaturalist 数据集效果显著。

本文的创新点如下:

(1) 引入基于 CLIP 的小样本 OOD 检测方法 LoCoOp, 通过抑制 ID 无关区域提升区分能力; 同时引入基于不确定性的自适应调制框架 SCT, 有效提高了分类准确率。

(2) 引入了 DPM 分数, 显著提高了在 iNaturalist 数据集上 ID 与 OOD 的区分性。

本文的局限性如下:

(1) 本文使用的模型局限于使用 ViT-B/16 作为视觉编码器的 CLIP 模型, 没有使用 ViT-B/32, ResNet-50 作为视觉编码器的 CLIP 模型, 也没有使用非对比式视觉语言模型如 BLIP 等, 对于方法的通用性没有进行充分的讨论。

(2) 本文没有详细讨论超参数的取值, 而这可能对结果造成较大影响, 同时也没有给出对于这些超参数的理论分析, 这不利于指导该方法的实际应用。

参考文献

- [1] Yang J, Zhou K, Li Y, et al. Generalized out-of-distribution detection: A survey[J]. International Journal of Computer Vision, 2024, 132(12): 5635-5662.
- [2] Miyai A, Yang J, Zhang J, et al. Generalized out-of-distribution detection and beyond in vision language model era: A survey[J]. arXiv preprint arXiv:2407.21794, 2024.
- [3] Lu S, Wang Y, Sheng L, et al. Recent Advances in OOD Detection: Problems and Approaches[J]. arXiv preprint arXiv:2409.11884, 2024.
- [4] Hendrycks D, Gimpel K. A Baseline for Detecting Misclassified and Out-of-Distribution Examples in Neural Networks[C]. International Conference on Learning Representations. 2017.
- [5] Hendrycks D, Basart S, Mazeika M, et al. Scaling Out-of-Distribution Detection for Real-World Settings[C]. International Conference on Machine Learning. PMLR, 2022: 8759-8773.
- [6] Liu W, Wang X, Owens J, et al. Energy-based out-of-distribution detection[J]. Advances in neural information processing systems, 2020, 33: 21464-21475.
- [7] Lee K, Lee K, Lee H, et al. A simple unified framework for detecting out-of-distribution samples and adversarial attacks[J]. Advances in neural information processing systems, 2018, 31.
- [8] Huang R, Geng A, Li Y. On the importance of gradients for detecting distributional shifts in the wild[J]. Advances in Neural Information Processing Systems, 2021, 34: 677-689.
- [9] Liang S, Li Y, Srikant R. Enhancing The Reliability of Out-of-distribution Image Detection in Neural Networks[C]. International Conference on Learning Representations. 2018.
- [10] Sun Y, Guo C, Li Y. React: Out-of-distribution detection with rectified activations[J]. Advances in neural information processing systems, 2021, 34: 144-157.
- [11] Wang H, Li Z, Feng L, et al. Vim: Out-of-distribution with virtual-logit matching[C]. Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2022: 4921-4930.
- [12] Wei H, Xie R, Cheng H, et al. Mitigating neural network overconfidence with logit normalization[C]. International conference on machine learning. PMLR, 2022: 23631-23644.
- [13] Zhang Z, Xiang X. Decoupling maxlogit for out-of-distribution detection[C]. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2023: 3388-3397.

- [14]Hsu Y C, Shen Y, Jin H, et al. Generalized odin: Detecting out-of-distribution image without learning from out-of-distribution data[C]. Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2020: 10951-10960.
- [15]Hendrycks D, Mazeika M, Dietterich T. Deep Anomaly Detection with Outlier Exposure[C]. International Conference on Learning Representations. 2019.
- [16]Zhang J, Inkawhich N, Linderman R, et al. Mixture outlier exposure: Towards out-of-distribution detection in fine-grained environments[C]. Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. 2023: 5531-5540.
- [17]Du X, Wang Z, Cai M, et al. VOS: Learning What You Don't Know by Virtual Outlier Synthesis[C]. International Conference on Learning Representations.2022.
- [18]Tao L, Du X, Zhu J, et al. Non-parametric Outlier Synthesis[C]. International Conference on Learning Representations.2023
- [19]Radford A, Kim J W, Hallacy C, et al. Learning transferable visual models from natural language supervision[C]. International conference on machine learning. PMLR, 2021: 8748-8763.
- [20]Ming Y, Cai Z, Gu J, et al. Delving into out-of-distribution detection with vision-language representations[J]. Advances in neural information processing systems, 2022, 35: 35087-35102.
- [21]Miyai A, Yu Q, Irie G, et al. GL-MCM: Global and Local Maximum Concept Matching for Zero-Shot Out-of-Distribution Detection[J]. International Journal of Computer Vision, 2025: 1-11.
- [22]Jiang X, Liu F, Fang Z, et al. Negative Label Guided OOD Detection with Pretrained Vision-Language Models[C]. The Twelfth International Conference on Learning Representations.
- [23]Zhou K, Yang J, Loy C C, et al. Learning to prompt for vision-language models[J]. International Journal of Computer Vision, 2022, 130(9): 2337-2348.
- [24]Ming Y, Li Y. How does fine-tuning impact out-of-distribution detection for vision-language models?[J]. International Journal of Computer Vision, 2024, 132(2): 596-609.
- [25]Miyai A, Yu Q, Irie G, et al. Locoop: Few-shot out-of-distribution detection via prompt learning[J]. Advances in Neural Information Processing Systems, 2024, 36.
- [26]Yu G, Zhu J, Yao J, et al. Self-Calibrated Tuning of Vision-Language Models for Out-of-Distribution Detection[J]. Advances in Neural Information Processing Systems, 2025, 37: 56322-56348.

- [27]Lafon M, Ramzi E, Rambour C, et al. Gallop: Learning global and local prompts for vision-language models[C]. European Conference on Computer Vision. Cham: Springer Nature Switzerland, 2024: 264-282.
- [28]Li T, Pang G, Bai X, et al. Learning transferable negative prompts for out-of-distribution detection[C]. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024: 17584-17594.
- [29]Bai Y, Han Z, Cao B, et al. Id-like prompt learning for few-shot out-of-distribution detection[C]. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024: 17480-17489.
- [30]Huang R, Li Y. Mos: Towards scaling out-of-distribution detection for large semantic space[C]. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021: 8710-8719.
- [31]Deng J, Dong W, Socher R, et al. Imagenet: A large-scale hierarchical image database[C]. 2009 IEEE conference on computer vision and pattern recognition. Ieee, 2009: 248-255.
- [32]Van Horn G, Mac Aodha O, Song Y, et al. The inaturalist species classification and detection dataset[C]. Proceedings of the IEEE conference on computer vision and pattern recognition. 2018: 8769-8778.
- [33]Xiao J, Hays J, Ehinger K A, et al. Sun database: Large-scale scene recognition from abbey to zoo[C]. 2010 IEEE computer society conference on computer vision and pattern recognition. IEEE, 2010: 3485-3492.
- [34]Zhou B, Lapedriza A, Khosla A, et al. Places: A 10 million image database for scene recognition[J]. IEEE transactions on pattern analysis and machine intelligence, 2017, 40(6): 1452-1464.
- [35]Cimpoi M, Maji S, Kokkinos I, et al. Describing textures in the wild[C]. Proceedings of the IEEE conference on computer vision and pattern recognition. 2014: 3606-3613.
- [36]Wang H, Li Y, Yao H, et al. Clipn for zero-shot ood detection: Teaching clip to say no[C]. Proceedings of the IEEE/CVF International Conference on Computer Vision. 2023: 1802-1812.

附录一 各方法在不同随机数种子下的具体结果

附表 1 各方法在 1-shot 设置下三个随机数种子的分布外检测性能结果

对比 方法	iNaturalist		SUN		Places		Texture		Average	
	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓
1-shot seed1										
CoOp	93.16	35.08	92.38	36.38	90.12	41.73	86.89	53.63	90.64	41.71
LoCoOp	94.06	29.57	94.59	25.56	91.42	35.02	87.16	53.07	91.81	35.80
SCT	93.77	30.91	94.18	27.46	91.33	35.24	86.07	54.34	91.34	36.99
DPCoOp	95.27	23.13	92.38	36.38	90.12	41.73	87.03	53.28	91.20	38.63
DPLoCoOp	96.12	18.47	94.59	25.56	91.42	35.02	87.50	52.36	92.41	32.85
DPSCT	96.22	17.55	94.18	27.46	91.33	35.24	86.47	53.26	92.05	33.38
1-shot seed2										
CoOp	92.12	40.09	92.96	35.07	90.05	43.21	87.73	50.02	90.71	42.10
LoCoOp	95.17	23.29	94.88	22.99	92.06	31.07	87.83	50.74	92.48	32.02
SCT	94.85	24.27	93.94	26.09	91.30	33.53	86.09	52.00	91.55	34.05
DPCoOp	92.58	39.55	92.96	35.07	90.05	43.21	87.73	50.02	90.83	41.96
DPLoCoOp	95.61	21.45	94.88	22.99	92.06	31.07	88.03	50.67	92.64	31.55
DPSCT	95.67	21.36	93.94	26.09	91.30	33.53	86.09	52.00	91.75	33.25
1-shot seed3										
CoOp	90.60	49.51	91.42	41.38	88.93	47.51	88.18	49.91	89.78	47.08
LoCoOp	95.14	24.30	94.63	24.27	91.71	33.62	87.69	50.62	92.29	33.20
SCT	95.64	20.41	93.89	27.24	91.37	34.56	87.69	49.68	92.15	32.97
DPCoOp	95.15	24.85	91.42	41.38	88.93	47.51	88.18	49.91	90.92	40.91
DPLoCoOp	97.17	13.22	94.63	24.27	91.71	33.62	87.69	50.62	92.80	30.43
DPSCT	97.61	10.55	93.89	27.24	91.37	34.56	87.93	49.75	92.70	30.53
1-shot average										
CoOp	91.96	41.56	92.25	37.61	89.70	44.15	87.60	51.19	90.38	43.63
LoCoOp	94.79	25.72	94.70	24.27	91.73	33.24	87.56	51.48	92.20	33.68
SCT	94.75	25.20	94.00	26.93	91.33	34.44	86.62	52.01	91.68	34.65
DpCoOp	94.33	29.18	92.25	37.61	89.70	44.15	87.65	51.07	90.98	40.50
DPLoCoOp	96.30	17.71	94.70	24.27	91.73	33.24	87.74	51.22	92.62	31.61
DPSCT	96.50	16.49	94.00	26.93	91.33	34.44	86.83	51.67	92.17	32.38

附表 2 各方法在 2-shot 设置下三个随机数种子的分布外检测性能结果

对比 方法	iNaturalist		SUN		Places		Texture		Average	
	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓
2-shot seed1										
CoOp	91.89	40.60	92.27	34.24	89.40	43.61	90.39	42.16	90.99	40.15
LoCoOp	95.86	18.88	94.39	26.16	91.79	34.01	89.72	43.53	92.94	30.64
SCT	95.84	19.72	94.09	25.42	91.11	34.34	88.36	45.46	92.35	31.24
DPCoOp	94.95	25.09	92.27	34.24	89.40	43.61	90.39	42.16	91.75	36.28
DPLoCoOp	97.26	12.30	94.39	26.16	91.79	34.01	89.92	43.72	93.34	29.05
DPSCT	97.43	11.59	94.09	25.42	91.11	34.34	88.63	45.04	92.81	29.10
2-shot seed2										
CoOp	93.88	29.71	91.73	36.98	89.44	43.85	88.03	49.73	90.77	40.07
LoCoOp	96.35	17.00	95.45	21.94	92.01	32.02	90.21	43.51	93.51	28.62
SCT	96.23	17.74	95.00	23.06	91.91	32.12	87.12	50.50	92.57	30.85
DPCoOp	95.51	21.99	91.73	36.98	89.44	43.85	88.14	49.15	91.20	37.99
DPLoCoOp	97.07	13.83	95.45	21.94	92.01	32.02	90.21	43.51	93.69	27.83
DPSCT	97.40	12.07	95.00	23.06	91.91	32.12	87.44	49.11	92.94	29.09
2-shot seed3										
CoOp	92.56	35.44	92.70	33.98	90.37	40.57	90.17	42.77	91.45	38.19
LoCoOp	97.00	14.05	95.50	20.60	92.74	30.12	89.98	44.20	93.81	27.24
SCT	94.89	24.71	94.62	24.72	91.67	33.12	87.75	50.28	92.23	33.21
DPCoOp	95.87	18.69	92.70	33.98	90.37	40.57	90.17	42.77	92.28	34.00
DPLoCoOp	97.62	10.95	95.50	20.60	92.74	30.12	90.26	42.96	94.03	26.16
DPSCT	97.09	12.71	94.62	24.72	91.67	33.12	88.09	49.80	92.87	30.09
2-shot average										
CoOp	92.78	35.25	92.23	35.07	89.74	42.68	89.53	44.89	91.07	39.47
LoCoOp	96.40	16.64	95.11	22.90	92.18	32.05	89.97	43.75	93.42	28.84
SCT	95.65	20.72	94.57	24.40	91.56	33.19	87.74	48.75	92.38	31.77
DPCoOp	95.44	21.92	92.23	35.07	89.74	42.68	89.57	44.69	91.75	36.09
DPLoCoOp	97.32	12.36	95.11	22.90	92.18	32.05	90.13	43.40	93.69	27.68
DPSCT	97.31	12.12	94.57	24.40	91.56	33.19	88.05	47.98	92.87	29.42

附表 3 各方法在 4-shot 设置下三个随机数种子的分布外检测性能结果

对比 方法	iNaturalist		SUN		Places		Texture		Average	
	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓
4-shot seed1										
CoOp	93.04	32.80	93.03	33.37	90.28	40.43	88.58	46.12	91.23	38.18
LoCoOp	96.44	16.11	94.90	24.50	91.75	33.86	89.80	42.87	93.22	29.34
SCT	96.03	17.61	95.71	19.64	92.74	28.17	88.62	44.93	93.28	27.59
DPCoOp	95.28	21.76	93.03	33.37	90.28	40.43	88.84	46.17	91.86	35.43
DPLoCoOp	97.32	11.81	94.90	24.50	91.75	33.86	90.01	42.62	93.50	28.20
DPSCT	97.39	11.01	95.71	19.64	92.74	28.17	88.89	44.73	93.68	25.89
4-shot seed2										
CoOp	93.53	30.46	92.23	36.29	89.70	43.29	88.84	45.98	91.08	39.00
LoCoOp	96.62	15.03	95.56	21.55	92.05	32.28	88.64	45.78	93.22	28.66
SCT	96.47	16.10	95.50	19.89	92.28	30.06	88.43	44.26	93.17	27.58
DPCoOp	95.27	22.47	92.23	36.29	89.70	43.29	88.97	45.60	91.54	36.91
DPLoCoOp	97.27	12.38	95.56	21.55	92.05	32.28	88.99	45.18	93.47	27.85
DPSCT	97.59	10.47	95.50	19.89	92.28	30.06	88.72	44.17	93.52	26.15
4-shot seed3										
CoOp	93.98	30.74	92.78	34.28	89.74	43.26	90.68	41.74	91.79	35.70
LoCoOp	95.53	20.95	95.52	20.77	92.53	30.17	89.67	43.03	93.31	28.73
SCT	96.49	16.41	94.77	23.04	91.61	33.00	89.45	43.16	93.08	28.90
DPCoOp	95.60	22.91	92.78	34.28	89.74	43.26	90.68	41.74	92.20	35.55
DPLoCoOp	96.80	15.11	95.52	20.77	92.53	30.17	89.83	42.98	93.67	27.26
DPSCT	97.69	10.06	94.77	23.04	91.61	33.00	89.45	43.16	93.38	27.31
4-shot average										
CoOp	93.52	31.33	92.68	34.65	89.91	42.33	89.37	44.61	91.37	38.23
LoCoOp	96.20	17.36	95.33	22.27	92.11	32.10	89.37	43.89	93.25	28.91
SCT	96.33	16.71	95.33	20.86	92.21	30.41	88.83	44.12	93.18	28.03
DPCoOp	95.38	22.38	92.68	34.65	89.91	42.33	89.50	44.50	91.87	35.97
DPLoCoOp	97.13	13.10	95.33	22.27	92.11	32.10	89.61	43.59	93.55	27.77
DPSCT	97.56	10.51	95.33	20.86	92.21	30.41	89.02	44.02	93.53	26.45

附表 4 各方法在 8-shot 设置下三个随机数种子的分布外检测性能结果

对比 方法	iNaturalist		SUN		Places		Texture		Average	
	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓
8-shot seed1										
CoOp	92.65	34.65	92.84	33.21	89.98	41.93	89.48	44.75	91.04	38.64
LoCoOp	94.92	23.19	95.68	20.62	92.95	29.12	89.79	45.14	93.33	29.52
SCT	95.22	21.55	94.95	21.56	92.39	30.63	89.39	44.38	92.99	29.53
DPCoOp	95.72	20.57	92.84	33.21	89.98	41.93	89.62	44.77	92.04	35.12
DPLoCoOp	97.01	13.12	95.68	20.62	92.95	29.12	90.00	44.86	93.91	26.93
DPSCT	97.16	12.55	94.95	21.56	92.39	30.63	89.62	43.72	93.53	27.12
8-shot seed2										
CoOp	94.20	29.88	93.84	28.67	91.12	37.14	89.55	46.26	92.18	35.49
LoCoOp	95.16	24.50	95.56	21.09	92.36	31.21	89.83	46.10	93.23	30.72
SCT	96.77	15.40	95.12	22.05	92.12	31.25	88.98	44.50	93.25	28.30
DPCoOp	95.73	20.93	93.84	28.67	91.12	37.14	89.55	46.26	92.56	33.25
DPLoCoOp	96.55	16.11	95.56	21.09	92.36	31.21	90.01	46.10	93.62	28.63
DPSCT	97.71	10.54	95.12	22.05	92.12	31.25	89.21	44.66	93.54	27.13
8-shot seed3										
CoOp	94.54	28.34	93.49	30.54	90.83	38.97	89.80	44.08	92.16	35.48
LoCoOp	97.12	13.77	95.00	23.09	91.95	33.31	89.73	44.18	93.45	28.59
SCT	96.50	16.82	94.82	23.38	92.29	31.73	89.28	42.98	93.22	28.73
DPCoOp	95.80	21.20	93.49	30.54	90.83	38.97	89.80	44.08	92.48	33.70
DPLoCoOp	97.55	11.33	95.00	23.09	91.95	33.31	89.94	44.04	93.61	27.94
DPSCT	97.48	11.23	94.82	23.38	92.29	31.73	89.53	42.13	93.53	27.12
8-shot average										
CoOp	93.80	30.96	93.39	30.81	90.64	39.35	89.61	45.03	91.86	36.54
LoCoOp	95.73	20.49	95.41	21.60	92.42	31.21	89.78	45.14	93.34	29.61
SCT	96.16	17.92	94.96	22.33	92.27	31.20	89.22	43.95	93.15	28.85
DPCoOp	95.75	20.90	93.39	30.81	90.64	39.35	89.66	45.04	92.36	34.03
DPLoCoOp	97.04	13.52	95.41	21.60	92.42	31.21	89.98	45.00	93.71	27.83
DPSCT	97.45	11.44	94.96	22.33	92.27	31.20	89.45	43.50	93.53	27.12

附表 5 各方法在 16-shot 设置下三个随机数种子的分布外检测性能结果

对比 方法	iNaturalist		SUN		Places		Texture		Average	
	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓
16-shot seed1										
CoOp	93.41	33.56	93.04	32.49	90.46	40.15	89.31	45.67	91.55	37.97
LoCoOp	95.53	20.48	94.86	23.34	92.48	30.42	90.54	39.79	93.55	28.51
SCT	95.53	19.97	94.89	23.07	92.37	29.90	89.03	43.28	92.95	29.06
DPCoOp	96.08	17.92	93.04	32.49	90.46	40.15	89.40	45.48	92.24	34.01
DPLoCoOp	97.14	12.64	94.86	23.34	92.48	30.42	90.54	39.79	93.75	26.55
DPSCT	97.37	11.07	94.89	23.07	92.37	29.90	89.03	43.28	93.42	26.83
16-shot seed2										
CoOp	95.10	23.97	92.56	34.16	89.60	43.53	90.06	44.17	91.83	36.46
LoCoOp	97.13	13.90	95.30	22.59	92.25	32.10	90.06	44.70	93.69	13.90
SCT	96.50	14.83	95.12	21.56	92.30	29.91	89.31	43.42	93.31	27.43
DpCoOp	96.49	16.79	92.56	34.16	89.60	43.53	90.11	43.76	92.19	34.56
DPLoCoOp	97.90	9.82	95.30	22.59	92.25	32.10	90.26	44.31	93.93	27.20
DPSCT	97.73	9.64	95.12	21.56	92.30	29.91	89.53	43.30	93.67	26.10
16-shot seed3										
CoOp	94.17	28.81	92.42	37.66	89.74	44.53	89.84	43.14	91.54	38.53
LoCoOp	96.67	16.04	94.54	23.89	91.75	34.21	90.67	39.91	93.41	28.51
SCT	96.94	14.02	95.64	19.92	92.63	29.67	89.78	41.91	93.75	26.38
DPCoOp	96.42	16.52	92.42	37.66	89.74	44.53	89.84	43.14	92.10	35.46
DPLoCoOp	97.64	10.52	94.54	23.89	91.75	34.21	90.84	39.70	93.69	27.08
DPSCT	98.07	8.69	95.64	19.92	92.63	29.67	89.78	41.91	94.03	25.05
16-shot average										
CoOp	94.23	28.78	92.67	34.77	89.93	42.74	89.74	44.33	91.64	37.66
LoCoOp	96.44	16.81	94.90	23.27	92.16	32.24	90.42	41.47	93.48	28.45
SCT	96.32	16.27	95.22	21.52	92.43	29.83	89.37	42.87	93.34	27.62
DPCoOp	96.33	17.08	92.67	34.77	89.93	42.74	89.78	44.13	92.18	34.68
DPLoCoOp	97.56	10.99	94.90	23.27	92.16	32.24	90.55	41.27	93.79	26.94
DPSCT	97.72	9.80	95.22	21.52	92.43	29.83	89.45	42.83	93.71	26.00

附录二 各方法在原始论文中报道的性能

附表 6 各方法在原始论文中报道的性能, 其中下标表示使用的 OOD 分数 (包括 MCM 和 GLMCM), CoOp 和 LoCoOp 的性能来自 LoCoOp 论文, SCT 的性能来自 SCT 论文

对比 方法	iNaturalist		SUN		Places		Texture		Average	
	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓	AUR↑	FPR↓
1-shot 微调										
CoOp _{MCM}	91.26	43.38	91.95	38.53	89.09	46.68	87.83	50.64	90.03	44.81
CoOp _{GLMCM}	95.27	21.30	92.16	31.66	89.31	40.44	84.25	52.93	90.25	36.58
LoCoOp _{MCM}	92.49	38.49	93.67	33.27	91.07	39.23	89.13	49.25	91.53	40.17
LoCoOp _{GLMCM}	94.89	24.61	94.59	25.62	92.12	34.00	87.49	49.86	92.14	33.52
SCT _{GLMCM}	95.70	19.16	94.58	23.52	91.23	32.81	86.66	48.87	92.04	31.09
2-shot 微调										
CoOp _{MCM}	92.12	38.89	91.58	39.38	88.98	44.18	89.16	44.92	90.46	41.85
CoOp _{GLMCM}	95.36	21.17	91.08	35.00	88.32	42.25	85.79	49.23	90.14	36.91
LoCoOp _{MCM}	92.76	35.38	93.31	33.95	90.38	41.15	89.76	45.07	91.55	38.89
LoCoOp _{GLMCM}	95.14	23.39	94.89	24.32	91.53	34.15	88.27	47.36	92.46	32.30
SCT _{GLMCM}	96.15	16.99	94.79	21.68	92.21	31.01	88.64	42.62	92.95	28.08
4-shot 微调										
CoOp _{MCM}	92.60	35.36	92.27	37.06	89.15	45.38	89.68	43.74	90.92	40.39
CoOp _{GLMCM}	95.52	18.95	92.90	29.58	89.64	38.72	85.87	48.03	90.98	33.82
LoCoOp _{MCM}	93.93	29.45	93.24	33.06	90.32	41.13	90.54	44.15	92.01	36.95
LoCoOp _{GLMCM}	96.07	18.49	95.00	22.85	91.86	32.38	89.10	44.72	93.01	29.61
SCT _{GLMCM}	97.03	13.88	94.85	22.13	92.09	30.20	87.96	45.53	92.98	27.93
8-shot 微调										
CoOp _{MCM}	92.96	35.17	92.50	34.45	89.76	41.17	89.92	43.29	91.29	38.52
CoOp _{GLMCM}	96.69	15.23	93.08	27.78	90.22	35.93	85.91	48.26	91.47	31.80
LoCoOp _{MCM}	94.60	27.12	93.23	33.87	90.53	40.53	90.98	42.49	92.34	36.00
LoCoOp _{GLMCM}	96.47	16.34	94.96	22.40	91.83	31.86	89.81	42.20	93.27	28.20
16-shot 微调										
CoOp _{MCM}	94.43	28.00	92.29	36.95	89.74	43.03	91.24	39.33	91.93	36.83
CoOp _{GLMCM}	96.62	14.60	92.65	28.48	89.98	36.49	88.03	43.13	91.82	30.67
LoCoOp _{MCM}	95.45	23.06	93.35	32.70	90.64	39.92	91.32	40.23	92.69	33.98
LoCoOp _{GLMCM}	96.86	16.05	95.07	23.44	91.98	32.87	90.19	42.28	93.52	28.66
SCT _{GLMCM}	95.86	13.94	95.33	20.55	92.24	29.86	89.06	41.51	93.37	26.47

研 究 生 签 字 梁一凡