
Delving into Out-of-Distribution Detection with Vision-Language Representations

Yifei Ming¹ Ziyang Cai¹ Jiuxiang Gu² Yiyou Sun¹ Wei Li³ Yixuan Li¹

¹Department of Computer Sciences, University of Wisconsin-Madison

²Adobe ³Google Research

{alvinming, ziyangc, sunyiyou, sharonli}@cs.wisc.edu

jigu@adobe.com mweili@google.com

Abstract

Recognizing out-of-distribution (OOD) samples is critical for machine learning systems deployed in the open world. The vast majority of OOD detection methods are driven by a single modality (*e.g.*, either vision or language), leaving the rich information in multi-modal representations untapped. Inspired by the recent success of vision-language pre-training, this paper enriches the landscape of OOD detection from a single-modal to a multi-modal regime. Particularly, we propose Maximum Concept Matching (MCM), a simple yet effective zero-shot OOD detection method based on aligning visual features with textual concepts. We contribute in-depth analysis and theoretical insights to understand the effectiveness of MCM. Extensive experiments demonstrate that MCM achieves superior performance on a wide variety of real-world tasks. MCM with vision-language features outperforms a common baseline with pure visual features on a hard OOD task with semantically similar classes by 13.1% (AUROC). Code is available at <https://github.com/deeplearning-wisc/MCM>.

1 Introduction

Out-of-distribution (OOD) detection is critical for deploying machine learning models in the wild, where samples from novel classes can naturally emerge and should be flagged for caution. Despite increasing attention, the vast majority of OOD detection methods are driven by single-modal learning [26, 29, 34, 68, 89, 93, 95, 98]. For example, labels are typically encoded as one-hot vectors in image classification, leaving the semantic information encapsulated in texts largely unexploited. OOD detection relying on pure visual information can inherit the limitations, *e.g.*, when an OOD input may be visually similar to in-distribution (ID) data yet semantically different from any ID class.

In this paper, we delve into a new landscape for OOD detection, departing from the classic single-modal toward a *multi-modal* regime. While the motivation is appealing, a core challenge remains: *how to effectively utilize joint vision-language features for OOD detection?* In the visual domain, existing methods typically require good feature representations [66, 72], and a distance metric under which OOD data points are relatively far away from the in-distribution (ID) data [42, 71]. These approaches, however, do not directly translate into the multi-modal regime. On the representation learning side, recent vision-language pre-training schemes such as CLIP [59] and ALIGN [33] have emerged as promising alternatives for visual representation learning. The main idea is to align an image with its corresponding textual description in the feature space. While the resulting representations are powerful, OOD detection based on such aligned multi-modal features is still in its infancy.

We bridge the gap by exploring a distance-based OOD detection approach, leveraging the joint vision-language representations. Our method capitalizes on the compatibility between visual features and textual features. By defining the textual features as the “*concept prototypes*” for each ID class,

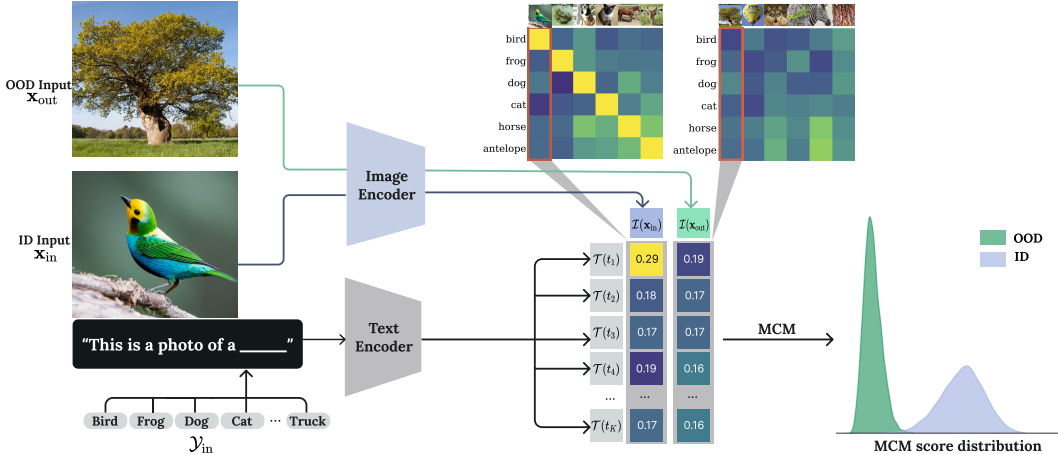


Figure 1: Overview of the proposed zero-shot OOD detection framework. The ID classification task is defined by a set of class labels \mathcal{Y}_{in} . The goal of OOD detection is to detect samples that do not belong to \mathcal{Y}_{in} . We view the textual embeddings of ID classes (wrapped by text templates) as concept prototypes. The OOD uncertainty of an input image can be characterized by the distance from visual features to the closest ID prototype. By properly scaling the distance, the MCM score achieves strong ID-OOD separability. See Section 3 for details.

we characterize OOD uncertainty by the distance from the visual feature to the closest ID prototype. That is, images closer to one of the textual embeddings of ID classes are more likely to be ID and vice versa. By a proper scaling of the distance, our proposed Maximum Concept Matching (MCM) score achieves strong ID-OOD separability (see Figure 1). MCM stands in contrast with the previous distance-based approaches, such as Mahalanobis [42], which defines class prototypes based on pure visual embeddings. Indeed, we show later in Section 5 that MCM (with multi-modal vision-language features) is far more competitive than Mahalanobis (with single-modal visual features). Moreover, while prior works of CLIP-based OOD detection [16, 19] rely on a set of candidate OOD labels, MCM is OOD-agnostic and alleviates the need for any prior information about test inputs.

Our work also advances the field by showcasing the promise of zero-shot OOD detection, which offers strong performance and generality without training on the ID samples. In particular, classic OOD detection methods often require training from scratch [9, 27] or fine-tuning [19, 32] on a given ID dataset. In this setting, a classifier and its companion OOD detector are good at only one task. Every new task (ID dataset) requires additional training and brings additional computation and storage costs. In contrast, we show for the first time that: (1) MCM achieves superior performance across a wide variety of real-world tasks—with just *one single pre-trained model*. This is encouraging given that there is no training or any OOD information involved. (2) On the challenging ImageNet-1k benchmark, MCM’s zero-shot OOD detection performance favorably matches and even outperforms strong task-specific baselines fine-tuned on BiT [32] and ViT models [19]. (3) MCM remains robust against hard OOD inputs, including both semantically hard OODs [85] and spurious OODs [50].

We summarize our main contributions as follows:

1. We propose MCM, a simple yet effective OOD detection method based on aligned vision-language features. MCM offers several compelling advantages over other OOD detection methods: generalizable (one model supports many tasks), OOD-agnostic (no information required from OOD data), training-free (no downstream fine-tuning required), and scalable to large real-world tasks.
2. We conduct extensive experiments and show that MCM achieves superior performance on a wide range of real-world tasks. On ImageNet-1k, MCM achieves an average AUROC of 91.49%, outperforming methods that require training. Moreover, we show that MCM remains competitive under challenging hard OOD evaluation tasks.
3. We provide in-depth empirical and theoretical analysis, providing insights to understand the effectiveness of MCM. We hope that this work will serve as a springboard for future works on OOD detection with multi-modal features.

2 Preliminaries

Contrastive vision-language pre-training. Compared to visual representation learning models such as ViT [13], vision-language representation learning demonstrates superior performance on image classification tasks. For instance, CLIP [59] adopts a self-supervised contrastive objective (*i.e.*, InfoNCE loss [75]) to align an image with its corresponding textual description in the feature space. Specifically, CLIP adopts a simple dual-stream architecture with one text encoder $\mathcal{T} : t \rightarrow \mathbb{R}^d$ (*e.g.*, Transformer [77]) and one image encoder $\mathcal{I} : \mathbf{x} \rightarrow \mathbb{R}^d$ (*e.g.*, ViT [13]). After pre-training on a dataset of 400 million text-image pairs, the joint vision-language embeddings of CLIP well associate objects in different modalities. Despite the promise, existing CLIP-like models perform zero-shot classification in a *closed-world* setting. That is, it will match an input into a fixed set of categories, even if it is irrelevant (*e.g.*, a tree being predicted as a bird in Figure 1). This motivates our work to leverage the multi-modal representation for OOD detection, which is largely unexplored.

Zero-shot OOD detection. Given a pre-trained model, a classification task of interest is defined by a set of class labels/names \mathcal{Y}_{in} , which we refer to as the known (ID) classes. Here ID classes are defined *w.r.t.* the classification task of interest, instead of the classes used in pre-training. Accordingly, OOD is defined *w.r.t.* the ID classes, not the data distribution during pre-training. The goal of OOD detection is to (1) detect samples that do not belong to any of the known classes; (2) otherwise, assign test samples to one of the known classes. Therefore, the OOD detector can be viewed as a “safeguard” for the classification model. Formally, we denote the OOD detector as a binary function: $G(\mathbf{x}; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) : \mathcal{X} \rightarrow \{\text{in}, \text{out}\}$, where $\mathbf{x} \in \mathcal{X}$ denotes a test image. Our method is based on only the names of the given classes in \mathcal{Y}_{in} , and a pre-trained model. Different from standard supervised learning, there is no training on the ID samples involved, hence zero-shot.

3 OOD Detection via Concept Matching

We illustrate our approach in Figure 1, which derives the OOD detector $G(\cdot)$ based on *concept matching*. For a given task with label set $\mathcal{Y}_{\text{in}} = \{y_1, y_2, \dots, y_K\}$, we can construct a collection of concept vectors $\mathcal{T}(t_i), i \in \{1, 2, \dots, K\}$, where t_i is the text prompt “this is a photo of a $\langle y_i \rangle$ ” for a label y_i . The concept vectors are represented by the embeddings of the text prompts.

For any test input image \mathbf{x}' , we can calculate the label-wise matching score based on the cosine similarity between the image feature $\mathcal{I}(\mathbf{x}')$ and the concept vector $\mathcal{T}(t_i)$: $s_i(\mathbf{x}') = \frac{\mathcal{I}(\mathbf{x}') \cdot \mathcal{T}(t_i)}{\|\mathcal{I}(\mathbf{x}')\| \cdot \|\mathcal{T}(t_i)\|}$. Formally, we define the maximum concept matching (**MCM**) score as:

$$S_{\text{MCM}}(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) = \max_i \frac{e^{s_i(\mathbf{x}')/\tau}}{\sum_{j=1}^K e^{s_j(\mathbf{x}')/\tau}}, \quad (1)$$

where τ is the temperature. For ID data, it will be matched to one of the concept vectors (textual prototypes) with a high score; and vice versa. Formally, our OOD detection function can be formulated as:

$$G(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) = \begin{cases} 1 & S_{\text{MCM}}(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) \geq \lambda \\ 0 & S_{\text{MCM}}(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) < \lambda \end{cases},$$

where by convention 1 represents the positive class (ID) and 0 indicates OOD. λ is chosen so that a high fraction of ID data (*e.g.*, 95%) is above the threshold. For samples that are classified as ID, one can obtain the class prediction based on the closest concept: $\hat{y} = \arg \max_{i \in [K]} s_i$.

Remark: (1) Our work differs from (and is complementary to) CLIP by focusing on OOD detection rather than (closed-world) zero-shot classification. We show new theoretical insights that softmax scaling plays a unique role in zero-shot OOD detection—improving the separability between ID and OOD data. This role has not been studied rigorously for zero-shot OOD detection. Readers familiar with CLIP may notice that MCM can be used for zero-shot classification in the closed world. This also makes MCM practically convenient for dual goals: detect OOD samples and assign ID data to one of the known classes. (2) Our method in principle is not limited to CLIP; it can be generally applicable for contrastive vision-language pre-training models that promote multi-modal feature alignment.

New insights on softmax scaling for zero-shot OOD detection. We provide theoretical justifications that softmax scaling improves the separability between ID and OOD data for CLIP-based OOD detection, which is *contrary* to models trained with cross-entropy (CE) loss. In particular, CLIP-like models are trained with a multi-modal contrastive loss, which maximizes the cosine similarity between an image and its textual description in the feature space. The resulting cosine similarity scores display strong *uniformity*¹ across labels, as evidenced in Figure 2 (right). Compared to OOD inputs, the gap between the maximum cosine similarity and the average is larger for ID inputs. However, the gap can be small when the number of ID classes increases where ID samples occur with lower highest cosine similarity. As a result, the highest cosine similarity for ID samples and OOD samples can be highly close (*c.f.* Figure 2 (left)).

Motivated by these observations, MCM employs softmax as a post hoc mechanism to **magnify** the difference. This is *fundamentally different from the softmax score derived from a model trained with cross-entropy loss*, which inherently maximizes the posterior $p(y|\mathbf{x})$ for the ground-truth label, and minimizes the probability for other labels. Unlike CLIP-like models, logit scores displaying uniformity would be heavily penalized by the CE loss. As a result, the logit score corresponding to the ground-truth label can already be significantly higher than other labels. Applying softmax on the logit scores can exacerbate overconfident predictions, and reduce the separability between ID and OOD data [46]. Indeed, for a model trained with cross-entropy loss, a logit-based score such as Energy [48] is shown to be much more effective than the softmax score.

Interestingly, for CLIP-like models, the trend is the opposite—applying softmax helps sharpen the uniform-like inner product scores, and increases the separability between ID and OOD data. To help readers better understand the insights, we first formalize our observations that OOD inputs trigger *similar cosine similarities* across ID concepts (Figure 2, right) as the following assumption:

Assumption 3.1. Let $z := \mathbb{1}\{y \in \mathcal{Y}_{\text{in}}\}$. $Q_{\mathbf{x}}$ denotes the out-of-distribution $\mathbb{P}_{\mathbf{x}|z=0}$ (marginal distribution of \mathbf{x} conditioned on $z = 0$). Assume $\exists \delta > 0$ such that

$$Q_{\mathbf{x}} \left(\frac{1}{K-1} \sum_{i \neq \hat{y}} [s_{\hat{y}_2}(\mathbf{x}) - s_i(\mathbf{x})] < \delta \right) = 1,$$

where $\hat{y} := \operatorname{argmax}_{i \in [K]} s_i(\mathbf{x})$ and $\hat{y}_2 := \operatorname{argmax}_{i \neq \hat{y}, i \in [K]} s_i(\mathbf{x})$ denote the indices of the largest and second largest cosine similarities for an OOD input \mathbf{x} .

Now we provide formal guarantees that using softmax can provably reduce the false positive rate (FPR) compared to that without softmax.

Theorem 3.1. Given a task with ID label set $\mathcal{Y}_{\text{in}} = \{y_1, y_2, \dots, y_K\}$ and a pre-trained CLIP-like model $(\mathcal{T}, \mathcal{I})$. If $Q_{\mathbf{x}}$ satisfies Assumption 3.1, then there exists a constant $T = \frac{\lambda(K-1)(\lambda^{\text{wo}} + \delta - s_{\hat{y}_2})}{K\lambda - 1}$ such that for any temperature $\tau > T$, we have

$$\text{FPR}(\tau, \lambda) \leq \text{FPR}^{\text{wo}}(\lambda^{\text{wo}}),$$

where $\text{FPR}(\tau, \lambda)$ is the false positive rate based on softmax scaling *with* temperature τ and detection threshold λ ; $\text{FPR}^{\text{wo}}(\lambda^{\text{wo}})$ is the false positive rate *without* softmax scaling based on threshold λ^{wo} .

¹This can be explained both theoretically [84] and empirically [81]. It has been shown that self-supervised contrastive learning with a smaller temperature (*e.g.*, initialized as 0.07 for CLIP) promotes uniform distribution for L_2 -normalized features. Moreover, as CLIP features lie on a high-dimensional space (512 for CLIP-B/16 and 768 for CLIP-L/14), uniformly distributed points in a high-dimensional sphere tend to be equidistant to each other [79]. Therefore, for OOD inputs, we observe approximately uniform cosine similarity with concept vectors.

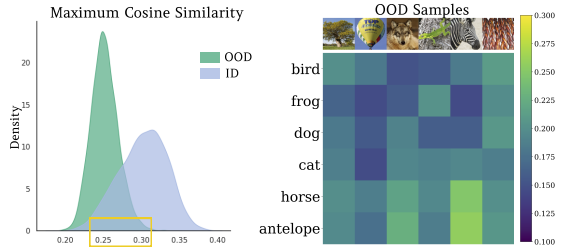


Figure 2: Left: Maximum cosine similarity for ID and OOD inputs. There exists overlapping regions (shown in yellow); Right: Cosine similarities between OOD inputs and ID concept vectors. For OOD inputs, the cosine similarities display uniformity.

This suggests that applying softmax scaling with a moderate temperature results in superior OOD detection performance compared to that without softmax scaling. The proof is in Appendix A. Later in Section 5, we empirically verify on a real-world ImageNet dataset that our bound can indeed be satisfied in CLIP where the thresholds are chosen at 95% true positive rate.

What MCM offers: Beyond theoretical insights, we would like to highlight several compelling advantages of our zero-shot OOD detection approach, owing to the strong pre-trained CLIP model:

- **Generalizable to many tasks:** Traditional OOD detection methods are based on a task-specific model. As a result, the OOD detector is not suitable for a realistic online scenario where the task changes from one to another. In contrast, we will show in Section 4 that MCM can perform a wide variety of OOD detection tasks, with just one single model. For a new task, only the names of the task’s visual concepts \mathcal{Y}_{in} are required.
- **OOD-agnostic:** Our method does not rely on any OOD information, and thus suits many real-world scenarios where one cannot anticipate what the unknowns would be ahead of time. This also mitigates the shortcoming of a recent approach [19], which assumes that a set of unseen labels are given as some weak information about OOD data.
- **Training-free:** MCM enables OOD detection in a zero-shot fashion. This stands in contrast to the vast majority of OOD detection literature, which often requires training from scratch or fine-tuning to achieve competitive performance.
- **Scalable:** The contrastive vision-language pre-training paradigm makes MCM scalable to a large number of class labels and realistic high-resolution images.

We now proceed to the experimental results, demonstrating these advantages on real-world tasks.

4 A Comprehensive Analysis of MCM

4.1 Datasets and Implementation Details

Datasets. Most previous works on OOD detection only focus on small-scale datasets with blurry images such as CIFAR [40] and TinyImageNet [41]. With pre-trained models such as CLIP, OOD detection can be extended to more realistic and complex datasets. In this work, we scale up evaluations in terms of (1) image resolution, (2) dataset variety, and (3) number of classes. We consider the following ID datasets: CUB-200 [80], STANFORD-CARS [39], FOOD-101 [6], OXFORD-PET [57] and variants of IMAGENET [11]. For OOD test datasets, we use the same ones in [32], including subsets of iNaturalist [76], SUN [86], PLACES [96], and TEXTURE [10]. For each OOD dataset, the categories are not overlapping with the ID dataset. We also use subsets of ImageNet-1k for fine-grained analysis. For example, we construct ImageNet-10 that mimics the class distribution of CIFAR-10 but with high-resolution images. For hard OOD evaluation, we curate ImageNet-20, which consists of 20 classes semantically similar to ImageNet-10 (*e.g.*, dog (ID) vs. wolf (OOD)).

Model. In our experiments, we adopt CLIP [59] as the target pre-trained model, which is one of the most popular and publicly available vision-language models. Note that our method is not limited to CLIP; it can generally be applicable for contrastive vision-language pre-training models that promote multi-modal feature alignment. Specifically, we mainly use CLIP-B/16, which consists of a ViT-B/16 Transformer as the image encoder and a masked self-attention Transformer [77] as the text encoder. To indicate the input patch size in ViT models, we append “/x” to model names. We prepend -B, -L to indicate Base and Large versions of the corresponding architecture. For instance, ViT-B/16 implies the Base variant with an input patch resolution of 16×16 . We also use CLIP-L/14 which is based on ViT-L/14 as a representative of large models. Unless specified otherwise, the temperature τ is 1 for all experiments. Details of the datasets, experimental setup, and hyperparameters are provided in Appendix B.

Metrics. For evaluation, we use the following metrics: (1) the false positive rate (FPR95) of OOD samples when the true positive rate of in-distribution samples is at 95%, (2) the area under the receiver operating characteristic curve (AUROC), and (3) ID classification accuracy (ID ACC).

Table 1: Zero-shot OOD detection with MCM score based on CLIP-B/16 with various ID datasets.

ID Dataset	OOD Dataset									
	iNaturalist		SUN		Places		Texture		Average	
	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑
CUB-200 [80]	9.83	98.24	4.93	99.10	6.65	98.57	6.97	98.75	7.09	98.66
Stanford-Cars [39]	0.05	99.77	0.02	99.95	0.24	99.89	0.02	99.96	0.08	99.89
Food-101 [6]	0.64	99.78	0.90	99.75	1.86	99.58	4.04	98.62	1.86	99.43
Oxford-Pet [57]	2.85	99.38	1.06	99.73	2.11	99.56	0.80	99.81	1.70	99.62
ImageNet-10	0.12	99.80	0.29	99.79	0.88	99.62	0.04	99.90	0.33	99.78
ImageNet-20	1.02	99.66	2.55	99.50	4.40	99.11	2.43	99.03	2.60	99.32
ImageNet-100	18.13	96.77	36.45	94.54	34.52	94.36	41.22	92.25	32.58	94.48

Table 2: OOD detection performance for ImageNet-1k [11] as ID.

Method	OOD Dataset									
	iNaturalist		SUN		Places		Texture		Average	
	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑
	Requires training (or w. fine-tuning)									
MOS [32] (BiT)	9.28	98.15	40.63	92.01	49.54	89.06	60.43	81.23	39.97	90.11
Fort et al. [19] (ViT-B)	15.07	96.64	54.12	86.37	57.99	85.24	53.32	84.77	45.12	88.25
Fort et al. [19] (ViT-L)	15.74	96.51	52.34	87.32	55.14	86.48	51.38	85.54	43.65	88.96
Energy [48] (CLIP-B)	21.59	95.99	34.28	93.15	36.64	91.82	51.18	88.09	35.92	92.26
Energy [48] (CLIP-L)	10.62	97.52	30.46	93.83	32.25	93.01	44.35	89.64	29.42	93.50
MSP [25] (CLIP-B)	40.89	88.63	65.81	81.24	67.90	80.14	64.96	78.16	59.89	82.04
MSP [25] (CLIP-L)	34.54	92.62	61.18	83.68	59.86	84.10	59.27	82.31	53.71	85.68
	Zero-shot (no training required)									
MCM (CLIP-B)	30.91	94.61	37.59	92.57	44.69	89.77	57.77	86.11	42.74	90.77
MCM (CLIP-L)	28.38	94.95	29.00	94.14	35.42	92.00	59.88	84.88	38.17	91.49

4.2 Main Results

MCM supports a diverse collection of tasks while being zero-shot. We first show that zero-shot OOD detection with MCM is effective across a wide variety of tasks—with just *one single pre-trained model*. To showcase the versatility of MCM, we consider the seven ID datasets here. To the best of our knowledge, this is among the first attempts to showcase the efficacy under an expansive and diverse collection of ID datasets. The zero-shot OOD detection performance is summarized in Table 1. A salient observation is that MCM can achieve superior detection performance on many tasks. For example, using STANFORD-CARS as ID, MCM yields an average FPR95 of **0.08%**. Considering that there are no training samples or OOD information involved, these results are very encouraging.

It can be also seen from Table 1 that MCM is promising, especially when the number of samples per ID class is limited in the training set. For example, there are only around 40 samples per class for Stanford-Cars, 100 for Oxford-Pet, and 30 for CUB-200. The sample scarcity makes OOD detection methods that rely on fine-tuning difficult. For example, after fine-tuning on Food-101, while the ID accuracy is increased from 86.3% to 92.5% ↑, OOD detection based on MSP is on par with MCM (99.5% vs. 99.4% in AUROC).

MCM scales effectively to large datasets. To examine the scalability of MCM, we compare it with recent competitive OOD detection methods [19, 32] on the ImageNet-1k dataset (ID) in Table 2. We observe the following trends:

- Larger models lead to superior performance. Compared with CLIP-B, MCM based on CLIP-L reduces FPR95 by 4.57%. Zero-shot ID classification accuracy is also improved by 6.27% with the larger model, reaching 73.28% (see Appendix D). This suggests that larger models are endowed with a better representation quality, which benefits both ID classification and OOD detection with MCM. Our finding echos with the recent observations [78] that higher ID classification accuracy is correlated with stronger OOD detection performance.
- MOS [32] recently demonstrated competitive performance on ImageNet-1k, which requires model fine-tuning based on BiT [38]. In contrast, we show that MCM (CLIP-L) outperforms MOS by 1.38% in AUROC while being zero-shot (training-free).
- MCM shares a softmax scaling function with the classic (visual) confidence-based score MSP [25]. To implement MSP, we adopt the commonly used linear probe approach by fine-tuning a linear layer on frozen visual features of CLIP. After fine-tuning, ID accuracy significantly improves, reaching 84.12% (CLIP-L). Interestingly, the OOD detection performance of MSP is worse than

Table 3: Performance comparison on **hard OOD detection** tasks. MCM is competitive on all three hard OOD tasks without training involved. MSP (based on fine-tuned CLIP) does not further improve performance.

Method	ID OOD	ImageNet-10	ImageNet-20	Waterbirds
		ImageNet-20	ImageNet-10	Spurious OOD
		FPR95 / AUROC	FPR95 / AUROC	FPR95 / AUROC
MSP [25] (fine-tuning)		9.38 / 98.31	12.51 / 97.70	39.57 / 90.99
Mahalanobis [42] (visual only)		78.32 / 85.60	43.03 / 89.94	2.21 / 99.55
MCM (zero-shot)		5.00 / 98.71	12.91 / 98.09	5.87 / 98.36

MCM by 15.54% in FPR95. Under the same model fine-tuned with linear probing, we observe that the Energy score outperforms MSP, corroborating findings in [48]. We investigate more in Section 5.

- Recently, Fort *et al.* [19] explore small-scale OOD detection by fine-tuning the full ViT model. When extended to large-scale tasks, we find that MCM still yields superior performance under the same image encoder configuration (ViT-B or ViT-L). This further highlights the advantage of utilizing vision-language joint embeddings for large-scale visual OOD detection.

MCM benefits hard OOD detection. Going beyond, we investigate whether MCM is still effective for hard OOD inputs. We consider the following two categories of hard OOD:

- **Semantically hard OOD:** OOD samples that are semantically similar to ID samples are particularly challenging for OOD detection algorithms [85]. To evaluate hard OOD detection tasks in realistic settings, here we consider ImageNet-10 (ID) vs. ImageNet-20 (OOD) and vice versa. The pair consists of high-resolution images with semantically similar categories such as dog versus wolf. As shown in Table 3, MCM outperforms Mahalanobis [42] by **73.32%** in FPR95 for ImageNet-10 (ID) vs. ImageNet-20 (OOD) and **30.12%** vice versa.
- **Spurious OOD:** Modern neural networks can exploit spurious correlations for predictions [3]. For example, in the Waterbirds dataset [64], there exist spurious correlations between the habitat (*e.g.*, water) and bird types. A recent work [50] proposes a new type of hard OOD named spurious OOD and shows that most OOD detection methods perform much worse for spurious OOD inputs compared to non-spurious inputs. The spurious OOD inputs are created to share the same background (*i.e.*, water) as ID data but have different object labels (*e.g.*, a boat rather than a bird). See Appendix C for illustrations. The results are shown in Table 3. It has been shown that CLIP representations are robust to distributional shifts [59]. Therefore, while prior works [50] show that spurious OOD inputs are challenging for methods based on ResNet [23], MCM and Mahalanobis scores based on pre-trained CLIP perform much better. On the other hand, fine-tuning exposes the model to the training set containing spurious correlations. As a result, MSP performs much worse than MCM (39.57% vs. 5.87% in FPR95).

MCM outperforms CLIP-based baselines. Two recent works also use CLIP embeddings for OOD detection [16, 19]. However, fundamental limitations exist for both works. Fort *et al.* [19] assume that a candidate OOD label set \mathcal{Y}_C is known, and used $\sum_{y \in \mathcal{Y}_C} \hat{p}(y|\mathbf{x})$ for OOD detection. Here the predictive probability $\hat{p}(y|\mathbf{x})$ is obtained by normalizing the inner products over $|\mathcal{Y}_{in}| + |\mathcal{Y}_C|$ classes. While applying softmax converts any vector to probabilities, as we show in Section 3, the converted probabilities do not necessarily correspond to $\mathbb{P}(\text{OOD}|\mathbf{x})$. Moreover, obtaining such an OOD label set is typically not feasible, which fundamentally limits its applicability. A recent work [16] realizes this idea by training an extra text decoder on top of CLIP’s image encoder to generate candidate labels. However, [16] cannot guarantee the generated labels are non-overlapping with the ID labels.

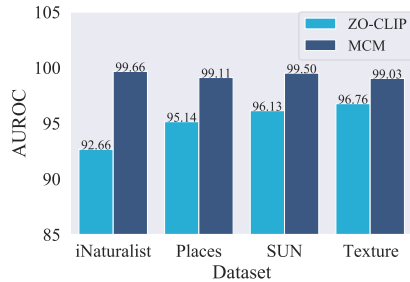


Figure 3: Comparison with a candidate label-based score ZO-CLIP on ImageNet-20, based on our implementation of [16]. Implementation details are deferred to Appendix E.1.

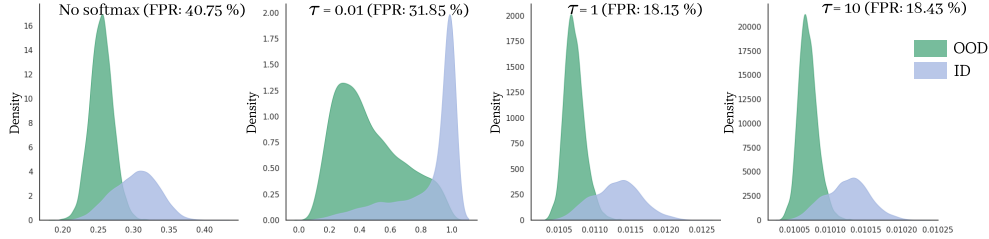


Figure 4: The influence of softmax scaling and temperature. We use ImageNet-100 (ID) vs. iNaturalist (OOD). Softmax scaling with a moderate temperature significantly improves FPR95.

We enhance the baseline with a stronger decoder and a filter module (see Appendix E.1). As shown in Figure 3, MCM outperforms the enhanced baseline on all OOD datasets. Moreover, MCM is much simpler to use—alleviating the need for an OOD label set or training an additional caption generator. In contrast, the caption generator’s performance largely affects OOD detection. Poor caption quality degenerates the OOD detection performance of candidate label-based methods. Moreover, obtaining a reliable caption generator for *any input image* can significantly increase the computational overhead.

5 Discussion: A Closer Look at MCM

Empirical verification on the role of softmax. In Section 3, we prove that softmax scaling on cosine similarity scores with a moderate τ improves the ID-OOD separability. Here we empirically verify our theoretical results. As shown in Figure 4, compared to directly using the maximum cosine similarity without softmax (leftmost figure), softmax scaling with a temperature $\tau = 1$ significantly improves the performance by 22.6% in FPR95, and further increasing τ (e.g., $\tau = 10$) leads to similar performance. The results are based on ImageNet-100 (ID) versus iNaturalist (OOD).

Now, we verify if our theoretical bound (c.f. Theorem 3.1) is satisfied empirically as well in Figure 4. From the leftmost figure, we can estimate $\lambda^{\text{wo}} \approx 0.26$, $\delta \approx 0.03$, and $s_{j_2} \approx 0.23$. By checking the third figure ($\tau = 1$ is the temperature value we use for most experiments), we approximate $\lambda \approx 0.011$. As $K = 100$, we plug in the values and obtain the lower bound $T = \frac{\lambda(K-1)(\lambda^{\text{wo}} + \delta - s_{j_2})}{K\lambda - 1} \approx 0.65$. Since $\tau = 1 > 0.65$, by Theorem 3.1, applying softmax scaling with $\tau = 1$ is provably superior to without softmax scaling for OOD detection.

Are vision-language features better than visual feature alone?

MCM can be interpreted as a distance-based approach—images that are closer to one of the K class prototypes are more likely to be ID and vice versa. Here the class prototypes are defined based on a textual encoder. Alternatively, one can define the class prototypes based on visual features. For example, Mahalanobis [42] defines a class prototype as the average of visual embeddings for images belonging to the same class. This raises the question whether MCM (with *multi-modal* vision-language features) is better than Mahalanobis (with *single-modal* visual feature). For a fair comparison, we use the same ViT image encoder from CLIP-B. Both MCM and Mahalanobis extract visual features from the penultimate layer. On ImageNet-1k, Mahalanobis displays a limited performance, with 73.14% AUROC averaged across four OOD test datasets (90.77% for MCM), as shown in Figure 5. From a practical perspective, Mahalanobis requires computing the inverse covariance matrix, which can be both computationally expensive and inaccurate when the number of samples is scarce and the number of ID classes grows. In contrast, MCM is easier to use and more robust.

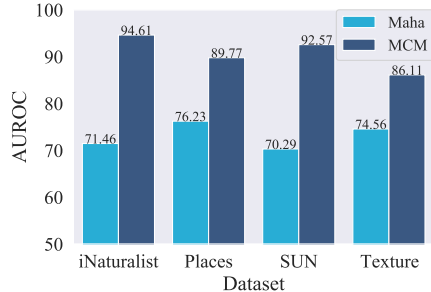


Figure 5: Comparison with Mahalanobis (Maha) score on ImageNet-1k.

MCM without softmax scaling. In Section 3, we provide theoretical justifications for the necessity of softmax scaling for CLIP-like models. To further verify our observations empirically, we show OOD detection performance based on the maximum cosine similarity score $S_{\text{MCM}}^{\text{wo}}(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) = \max_{i \in [K]} s_i(\mathbf{x}')$. The results are shown in Table 4. For easy tasks such as Food-101 [39], Stanford-

Table 4: Zero-shot OOD detection of S_{MCM}^{wo} based on CLIP-B/16.

ID Dataset	OOD Dataset									
	iNaturalist		SUN		Places		Texture		Average	
	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑
Stanford-Cars [39]	0.00	100	0.02	99.99	0.26	99.94	0.00	100	0.07	99.98
Food-101 [6]	0.56	99.86	0.09	99.95	0.49	99.88	8.33	97.44	2.37	99.28
Oxford-Pet [57]	0.02	99.98	0.05	99.97	0.20	99.94	0.27	99.91	0.14	99.95
ImageNet-10	2.40	99.42	1.79	99.55	2.83	99.32	1.86	99.56	2.22	99.46
ImageNet-20	14.96	97.87	13.10	97.97	14.21	97.67	13.46	97.32	13.93	97.71
ImageNet-1k	61.66	89.31	64.39	87.43	63.67	85.95	86.61	71.68	69.08	83.59

Cars [39], and Oxford-Pet [57] as ID, the performance of maximum cosine similarity score is similar to MCM (see Table 1 and Table 2). However, for more challenging tasks such as ImageNet-20 and ImageNet-1k, MCM significantly outperforms that without softmax scaling. For example, the average FPR95 is improved by **11.33%** on ImageNet-20 and **26.34%** on ImageNet-1k, which highlights the necessity of a proper scaling function for CLIP-based OOD detection.

MCM for ResNet-based CLIP models. Our main results are based on the CLIP model with ViT image encoder. We additionally investigate the effectiveness of MCM on ResNet-based CLIP. Specifically, we use RN50x4 (178.3M), which shares a similar number of parameters as CLIP-B/16 (149.6M). The results are shown in Table 5. We can see that MCM still shows promising results with ResNet-based CLIP models, and the performance is comparable between RN50x4 and CLIP-B/16 (89.97 vs. 90.77 in AUROC).

Table 5: Comparison with ResNet-based CLIP models on ImageNet-1k (ID).

Model	OOD Dataset									
	iNaturalist		SUN		Places		Texture		Average	
	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑
RN50x4	44.51	91.51	35.11	92.84	43.74	89.60	57.73	85.93	45.27	89.97
CLIP-B/16	30.91	94.61	37.59	92.57	44.69	89.77	57.77	86.11	42.74	90.77

Effect of prompt ensembling. We examine MCM’s performance with prompt ensembling. For example, Radford *et al.* [59] create 80 possible prompts according to the image modalities and nuances in ImageNet. We experiment with the two prompt sets, one of size 80 as in [59], and our own set of 5 prompts. Ensembles are obtained by averaging the textual features. As expected, using ensembles increases the ID classification accuracy on ImageNet-1k (2% with CLIP-B and 3% with CLIP-L). For OOD detection, the average FPR95 is reduced from 38.17% with the default prompt to 35.23%↓ with an ensemble of five prompts shown in Table 6. In addition, the detection performance with 5 prompts is slightly better than with 80 prompts. Note that prompt ensembling does not increase the inference-time cost, as the textual embeddings (across many prompts) can be pre-calculated and averaged into a single embedding.

A photo of a <label>.
A blurry photo of a <label>.
A photo of many <label>.
A photo of the large <label>.
A photo of the small <label>.

Table 6: The five prompt templates.

6 Related Works

OOD detection in computer vision. For open-world multi-class classification, the goal of OOD detection is to derive a binary ID-OOD classifier along with a multi-class classification model for visual inputs. A plethora of methods has been proposed for deep neural networks [91], including generative model-based methods [7, 20, 36, 53, 54, 56, 61, 67, 88], and discriminative-model based methods. For the latter category, an OOD score can be derived based on the softmax output [4, 12, 24, 25, 29, 32, 46, 90], energy-based score [15, 48, 49, 69, 70, 82], gradient information [31], or the feature embeddings [14, 42, 65, 66, 71, 72, 85] of a model. Morteza *et al.* [52], Fang *et al.* [17], and Bitterwolf *et al.* [5] provided theoretical analysis for OOD detection. Recent works [63, 83] also explored OOD detection for long-tailed distributions. Works insofar have mostly focused on OOD detection for a task-specific model using only visual information. In contrast, we explore a novel

paradigm of zero-shot OOD detection that incorporates rich textual information and can perform a wide variety of tasks.

OOD detection in natural language processing. Distribution shifts can occur due to the change of topics and domains, unexpected user utterances, *etc.* Challenging benchmarks [37] and characterization of distributional shifts [1] have been proposed in recent years. Compared to early language models such as ConvNets and LSTM [28], pre-trained language models are more robust to distribution shifts and more effective at identifying OOD instances [26, 58, 89]. Various algorithmic solutions are proposed to handle OOD detection, including outlier exposure [30], model ensembling [44], data augmentation [8, 93, 95], contrastive learning [34, 98], and an auxiliary module that incorporates domain labels [68]. Tan *et al.* [73] also explore zero-shot OOD detection for text classification tasks. However, prior works focus on pure natural language processing (NLP) settings, while we explore utilizing textual embeddings for zero-shot *visual* OOD detection.

Vision-language models. Utilizing large-scale pre-trained vision-language models for multimodal downstream tasks has become an emerging paradigm with remarkable performance [22, 74]. In general, two types of architectures exist: single-stream models like VisualBERT [43] and ViLT [35] feed the concatenated text and visual features into a single transformer-based encoder; dual-stream models such as CLIP [59], ALIGN [33], and FILIP [92] use separate encoders for text and image and optimize with contrastive objectives to align semantically similar features in different modalities. In particular, CLIP enjoys popularity due to its simplicity and strong performance. CLIP-like models inspire numerous follow-up works [45, 94, 97], which aim to improve data efficiency and better adaptation to downstream tasks. This paper adopts CLIP as the target pre-trained model, but our approach can be generally applicable to contrastive models that promote vision-language alignment.

Multi-modal OOD detection. Exploring textual information for visual OOD detection is a new area with limited existing works. Fort *et al.* [19] propose to feed the potential OOD labels to the textual encoder of CLIP [59]. Recently, Esmaeilpour *et al.* [16] propose to train a label generator based on the visual encoder of CLIP and use the generated labels for OOD detection. While both works rely on a set of candidate OOD labels, MCM is OOD-agnostic and alleviates the need for prior information on OOD. Moreover, prior works [16, 59] only focus on small-scale inputs. We largely expand the scope to a wide range of large-scale realistic datasets, and show new theoretical insights.

7 Conclusion

In this work, we delve into a new landscape for OOD detection, departing from the classic single-modal toward a multi-modal regime. By viewing the textual features as the “concept prototypes”, we explore a new OOD detection approach MCM, based on the joint vision-language representations. Unlike the majority of OOD detection methods, MCM offers several compelling advantages: training-free, generalizable to many tasks, scalable to hundreds of classes, and does not require any prior information on OOD inputs. Moreover, we provide theoretical guarantees on how softmax scaling provably improves zero-shot OOD detection. We investigate the effectiveness of MCM on a wide range of large-scale realistic tasks, including several types of hard OOD datasets. Lastly, we demonstrate the advantage of vision-language features over pure visual features for OOD detection. We hope our work will inspire future research toward multi-modal OOD detection.

Acknowledgement

The authors wish to thank Junjie Hu, Ying Fan, Ruisu Zhang, Andrew Geng, and Soumya Suvra Ghosal for the helpful discussions. The work is supported by a Google-Initiated Research Grant, and gift funding from Adobe Research.

References

- [1] Udit Arora, William Huang, and He He. Types of out-of-distribution texts and how to detect them. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2021.
- [2] Andrei Barbu, David Mayo, Julian Alverio, William Luo, Christopher Wang, Dan Gutfreund, Josh Tenenbaum, and Boris Katz. Objectnet: A large-scale bias-controlled dataset for pushing

- the limits of object recognition models. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2019.
- [3] Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In *The European Conference on Computer Vision (ECCV)*, 2018.
 - [4] Abhijit Bendale and Terrance E Boult. Towards open set deep networks. In *The IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*, 2016.
 - [5] Julian Bitterwolf, Alexander Meinke, Maximilian Augustin, and Matthias Hein. Breaking down out-of-distribution detection: Many methods based on ood training data estimate a combination of the same core quantities. In *International Conference on Machine Learning*, pages 2041–2074. PMLR, 2022.
 - [6] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative components with random forests. In *The European Conference on Computer Vision (ECCV)*, 2014.
 - [7] Mu Cai and Yixuan Li. Out-of-distribution detection via frequency-regularized generative models. In *Proceedings of IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023.
 - [8] Derek Chen and Zhou Yu. Gold: improving out-of-scope detection in dialogues using data augmentation. *arXiv preprint arXiv:2109.03079*, 2021.
 - [9] Jiefeng Chen, Yixuan Li, Xi Wu, Yingyu Liang, and Somesh Jha. Atom: Robustifying out-of-distribution detection using outlier mining. In *The European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PKDD)*, 2021.
 - [10] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *The IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*, 2014.
 - [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *The IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*, 2009.
 - [12] Terrance DeVries and Graham W Taylor. Learning confidence for out-of-distribution detection in neural networks. *arXiv preprint arXiv:1802.04865*, 2018.
 - [13] 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. In *International Conference on Learning Representations (ICLR)*, 2021.
 - [14] Xuefeng Du, Gabriel Gozum, Yifei Ming, and Yixuan Li. Siren: Shaping representations for detecting out-of-distribution objects. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
 - [15] Xuefeng Du, Zhaoning Wang, Mu Cai, and Yixuan Li. Vos: Learning what you don’t know by virtual outlier synthesis. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2022.
 - [16] Sepideh Esmailpour, Bing Liu, Eric Robertson, and Lei Shu. Zero-shot open set detection by extending clip. In *The AAAI Conference on Artificial Intelligence (AAAI)*, 2022.
 - [17] Zhen Fang, Yixuan Li, Jie Lu, Jiahua Dong, Bo Han, and Feng Liu. Is out-of-distribution detection learnable? In *Advances in Neural Information Processing System (NeurIPS)*, 2022.
 - [18] Christiane Fellbaum. Wordnet. In *Theory and Applications of Ontology: Computer Applications*, pages 231–243. Springer, 2010.
 - [19] Stanislav Fort, Jie Ren, and Balaji Lakshminarayanan. Exploring the limits of out-of-distribution detection. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2021.
 - [20] ZongYuan Ge, Sergey Demyanov, Zetao Chen, and Rahil Garnavi. Generative openmax for multi-class open set classification. *arXiv preprint arXiv:1707.07418*, 2017.
 - [21] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and Wieland Brendel. Imagenet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations (ICLR)*, 2019.

- [22] Jiuxiang Gu, Jason Kuen, Shafiq Joty, Jianfei Cai, Vlad Morariu, Handong Zhao, and Tong Sun. Self-supervised relationship probing. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *The IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*, 2016.
- [24] Matthias Hein, Maksym Andriushchenko, and Julian Bitterwolf. Why relu networks yield high-confidence predictions far away from the training data and how to mitigate the problem. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, pages 41–50, 2019.
- [25] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *International Conference on Learning Representations (ICLR)*, 2017.
- [26] Dan Hendrycks, Xiaoyuan Liu, Eric Wallace, Adam Dziedzic, Rishabh Krishnan, and Dawn Song. Pretrained transformers improve out-of-distribution robustness. In *Association for Computational Linguistics (ACL)*, 2020.
- [27] Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier exposure. In *International Conference on Learning Representations (ICLR)*, 2018.
- [28] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [29] Yen-Chang Hsu, Yilin Shen, Hongxia Jin, and Zsolt Kira. Generalized odin: Detecting out-of-distribution image without learning from out-of-distribution data. In *The IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*, 2020.
- [30] Yibo Hu and Latifur Khan. Uncertainty-aware reliable text classification. In *SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, 2021.
- [31] Rui Huang, Andrew Geng, and Yixuan Li. On the importance of gradients for detecting distributional shifts in the wild. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2021.
- [32] Rui Huang and Yixuan Li. Mos: Towards scaling out-of-distribution detection for large semantic space. In *The IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*, 2021.
- [33] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning (ICML)*, 2021.
- [34] Di Jin, Shuyang Gao, Seokhwan Kim, Yang Liu, and Dilek Hakkani-Tur. Towards textual out-of-domain detection without in-domain labels. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2022.
- [35] Wonjae Kim, Bokyoung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without convolution or region supervision. In *International Conference on Machine Learning (ICML)*, 2021.
- [36] Polina Kirichenko, Pavel Izmailov, and Andrew G Wilson. Why normalizing flows fail to detect out-of-distribution data. *Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [37] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning (ICML)*, 2021.
- [38] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. In *The European Conference on Computer Vision (ECCV)*, 2020.
- [39] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *4th International IEEE Workshop on 3D Representation and Recognition (3dRR-13)*, Sydney, Australia, 2013.
- [40] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.

- [41] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- [42] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2018.
- [43] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- [44] Xiaoya Li, Jiwei Li, Xiaofei Sun, Chun Fan, Tianwei Zhang, Fei Wu, Yuxian Meng, and Jun Zhang. kfolden: k-fold ensemble for out-of-distribution detection. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2021.
- [45] Yangguang Li, Feng Liang, Lichen Zhao, Yufeng Cui, Wanli Ouyang, Jing Shao, Fengwei Yu, and Junjie Yan. Supervision exists everywhere: A data efficient contrastive language-image pre-training paradigm. In *International Conference on Learning Representations (ICLR)*, 2022.
- [46] Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. In *International Conference on Learning Representations (ICLR)*, 2018.
- [47] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *The European Conference on Computer Vision (ECCV)*, 2014.
- [48] Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [49] Yifei Ming, Ying Fan, and Yixuan Li. Poem: Out-of-distribution detection with posterior sampling. In *International Conference on Machine Learning (ICML)*, 2022.
- [50] Yifei Ming, Hang Yin, and Yixuan Li. On the impact of spurious correlation for out-of-distribution detection. *The AAAI Conference on Artificial Intelligence (AAAI)*, 2022.
- [51] Ron Mokady, Amir Hertz, and Amit H Bermano. Clipcap: Clip prefix for image captioning. *arXiv preprint arXiv:2111.09734*, 2021.
- [52] Peyman Morteza and Yixuan Li. Provable guarantees for understanding out-of-distribution detection. *The AAAI Conference on Artificial Intelligence (AAAI)*, 2022.
- [53] Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, and Balaji Lakshminarayanan. Do deep generative models know what they don't know? In *International Conference on Learning Representations (ICLR)*, 2019.
- [54] Lawrence Neal, Matthew Olson, Xiaoli Fern, Weng-Keen Wong, and Fuxin Li. Open set learning with counterfactual images. In *The European Conference on Computer Vision (ECCV)*, 2018.
- [55] Edwin G. Ng, Bo Pang, Piyush Sharma, and Radu Soricut. Understanding guided image captioning performance across domains. *arXiv preprint arXiv:2012.02339*, 2020.
- [56] Poojan Oza and Vishal M Patel. C2ae: Class conditioned auto-encoder for open-set recognition. In *The IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*, 2019.
- [57] Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Cats and dogs. In *The IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*, 2012.
- [58] Alexander Podolskiy, Dmitry Lipin, Andrey Bout, Ekaterina Artemova, and Irina Piontkovskaya. Revisiting mahalanobis distance for transformer-based out-of-domain detection. In *The AAAI Conference on Artificial Intelligence (AAAI)*, 2021.
- [59] 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 (ICML)*, 2021.
- [60] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [61] Jie Ren, Peter J Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark Depristo, Joshua Dillon, and Balaji Lakshminarayanan. Likelihood ratios for out-of-distribution detection. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2019.

- [62] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?" explaining the predictions of any classifier. In *SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, 2016.
- [63] Abhijit Guha Roy, Jie Ren, Shekoofeh Azizi, Aaron Loh, Vivek Natarajan, Basil Mustafa, Nick Pawlowski, Jan Freyberg, Yuan Liu, Zach Beaver, et al. Does your dermatology classifier know what it doesn't know? detecting the long-tail of unseen conditions. *Medical Image Analysis*, 75:102274, 2022.
- [64] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. In *International Conference on Learning Representations (ICLR)*, 2019.
- [65] Chandramouli Shama Sastry and Sageev Oore. Detecting out-of-distribution examples with Gram matrices. In *International Conference on Machine Learning (ICML)*, 2020.
- [66] Vikash Sehwal, Mung Chiang, and Prateek Mittal. Ssd: A unified framework for self-supervised outlier detection. In *International Conference on Learning Representations (ICLR)*, 2021.
- [67] Joan Serra, David Álvarez, Vicenç Gómez, Olga Slizovskaia, José F. Núñez, and Jordi Luque. Input complexity and out-of-distribution detection with likelihood-based generative models. In *International Conference on Learning Representations (ICLR)*, 2020.
- [68] Yilin Shen, Yen-Chang Hsu, Avik Ray, and Hongxia Jin. Enhancing the generalization for intent classification and out-of-domain detection in slu. *arXiv preprint arXiv:2106.14464*, 2021.
- [69] Yiyu Sun, Chuan Guo, and Yixuan Li. React: Out-of-distribution detection with rectified activations. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2021.
- [70] Yiyu Sun and Yixuan Li. Dice: Leveraging sparsification for out-of-distribution detection. In *Proceedings of European Conference on Computer Vision (ECCV)*, 2022.
- [71] Yiyu Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. Out-of-distribution detection with deep nearest neighbors. In *International Conference on Machine Learning (ICML)*, 2022.
- [72] Jihoon Tack, Sangwoo Mo, Jongheon Jeong, and Jinwoo Shin. Csi: Novelty detection via contrastive learning on distributionally shifted instances. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [73] Ming Tan, Yang Yu, Haoyu Wang, Dakuo Wang, Saloni Potdar, Shiyu Chang, and Mo Yu. Out-of-domain detection for low-resource text classification tasks. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2019.
- [74] Shagun Uppal, Sarthak Bhagat, Devamanyu Hazarika, Navonil Majumder, Soujanya Poria, Roger Zimmermann, and Amir Zadeh. Multimodal research in vision and language: A review of current and emerging trends. *Information Fusion*, 77:149–171, 2022.
- [75] Aaron Van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv e-prints*, pages arXiv–1807, 2018.
- [76] 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 *The IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*, 2018.
- [77] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Conference on Neural Information Processing Systems (NeurIPS)*, 2017.
- [78] 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 (ICLR)*, 2022.
- [79] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- [80] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The caltech-ucsd birds-200-2011 dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011.
- [81] Feng Wang and Huaping Liu. Understanding the behaviour of contrastive loss. In *The IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*, 2021.

- [82] Haoran Wang, Weitang Liu, Alex Bocchieri, and Yixuan Li. Can multi-label classification networks know what they don't know? *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [83] Haotao Wang, Aston Zhang, Yi Zhu, Shuai Zheng, Mu Li, Alex J Smola, and Zhangyang Wang. Partial and asymmetric contrastive learning for out-of-distribution detection in long-tailed recognition. In *International Conference on Machine Learning (ICML)*, 2022.
- [84] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning (ICML)*, 2020.
- [85] Jim Winkens, Rudy Bunel, Abhijit Guha Roy, Robert Stanforth, Vivek Natarajan, Joseph R Led-sam, Patricia MacWilliams, Pushmeet Kohli, Alan Karthikesalingam, Simon Kohl, et al. Contrastive training for improved out-of-distribution detection. *arXiv preprint arXiv:2007.05566*, 2020.
- [86] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *The IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*, 2010.
- [87] Kai Yuanqing Xiao, Logan Engstrom, Andrew Ilyas, and Aleksander Madry. Noise or signal: The role of image backgrounds in object recognition. In *International Conference on Learning Representations (ICLR)*, 2021.
- [88] Zhisheng Xiao, Qing Yan, and Yali Amit. Likelihood regret: An out-of-distribution detection score for variational auto-encoder. In *Conference on Neural Information Processing Systems (NeurIPS)*, volume 33, 2020.
- [89] Keyang Xu, Tongzheng Ren, Shikun Zhang, Yihao Feng, and Caiming Xiong. Unsupervised out-of-domain detection via pre-trained transformers. In *Association for Computational Linguistics (ACL)*, 2021.
- [90] Jingkang Yang, Haoqi Wang, Litong Feng, Xiaopeng Yan, Huabin Zheng, Wayne Zhang, and Ziwei Liu. Semantically coherent out-of-distribution detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021.
- [91] Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-distribution detection: A survey. *arXiv preprint arXiv:2110.11334*, 2021.
- [92] Lewei Yao, Runhui Huang, Lu Hou, Guansong Lu, Minzhe Niu, Hang Xu, Xiaodan Liang, Zhenguo Li, Xin Jiang, and Chunjing Xu. Filip: Fine-grained interactive language-image pre-training. *International Conference on Learning Representations (ICLR)*, 2021.
- [93] Li-Ming Zhan, Haowen Liang, Bo Liu, Lu Fan, Xiao-Ming Wu, and Albert Lam. Out-of-scope intent detection with self-supervision and discriminative training. *Association for Computational Linguistics (ACL)*, 2021.
- [94] Renrui Zhang, Rongyao Fang, Peng Gao, Wei Zhang, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li. Tip-adapter: Training-free clip-adapter for better vision-language modeling. *arXiv preprint arXiv:2111.03930*, 2021.
- [95] Yinhe Zheng, Guanyi Chen, and Minlie Huang. Out-of-domain detection for natural language understanding in dialog systems. *TASLP*, 2020.
- [96] 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 (TPAMI)*, 2017.
- [97] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *The IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*, 2022.
- [98] Wenxuan Zhou, Fangyu Liu, and Muhao Chen. Contrastive out-of-distribution detection for pretrained transformers. *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2021.
- [99] Zhuotun Zhu, Lingxi Xie, and Alan Yuille. Object recognition with and without objects. In *International Joint Conferences on Artificial Intelligence (IJCAI)*, 2017.

A Theoretical Justification: Softmax Scaling for Zero-Shot OOD Detection

In this section, we provide the proof for Theorem 3.1 in Section 3, which states the benefits of applying softmax scaling to inner products for OOD detection. We begin with a review of notations.

Notations. We denote the text encoder of a pre-trained CLIP-like model as $\mathcal{T} : t \rightarrow \mathbb{R}^d$ and the image encoder $\mathcal{I} : \mathbf{x} \rightarrow \mathbb{R}^d$. For a given task with label set $\mathcal{Y}_{\text{in}} = \{y_1, y_2, \dots, y_K\}$, we construct a collection of concept vectors $\mathcal{T}(t_i)$. For a given input \mathbf{x}' , we denote the cosine similarity *w.r.t.* concept vectors as $s_i(\mathbf{x}') = \frac{\mathcal{I}(\mathbf{x}') \cdot \mathcal{T}(t_i)}{\|\mathcal{I}(\mathbf{x}')\| \cdot \|\mathcal{T}(t_i)\|} \forall i \in [K]$, where $|s_i(\mathbf{x}')| \leq B$ for all $\mathbf{x}' \in \mathcal{X}$.² We define the maximum concept matching (MCM) score as: $S_{\text{MCM}}(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) = \max_{i \in [K]} \frac{e^{s_i(\mathbf{x}')/\tau}}{\sum_{j=1}^K e^{s_j(\mathbf{x}')/\tau}}$. We denote the maximum inner product without applying softmax scaling as $S_{\text{MCM}}^{\text{wo}}(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) = \max_{i \in [K]} s_i(\mathbf{x}')$. By convention, the OOD detection functions are given by:

$$G^{\text{wo}}(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) = \begin{cases} 1 & S_{\text{MCM}}^{\text{wo}}(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) \geq \lambda^{\text{wo}} \\ 0 & S_{\text{MCM}}^{\text{wo}}(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) < \lambda^{\text{wo}} \end{cases},$$

$$G(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) = \begin{cases} 1 & S_{\text{MCM}}(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) \geq \lambda \\ 0 & S_{\text{MCM}}(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) < \lambda \end{cases},$$

Remarks: By convention, 1 represents the positive class (ID) and 0 indicates OOD; λ and λ^{wo} are typically chosen such that the true positive rate is at 95%.

For convenience, we paste the assumptions and the theorem in Section 3 below,

Assumption A.1. Let $z := \mathbb{1}\{y \in \mathcal{Y}_{\text{in}}\}$ and $Q_{\mathbf{x}}$ denotes the out-of-distribution $\mathbb{P}_{\mathbf{x}|z=0}$ (marginal distribution of \mathbf{x} conditioned on $z = 0$). Assume $\exists \delta > 0$ such that

$$Q_{\mathbf{x}} \left(\frac{1}{K-1} \sum_{i \neq \hat{y}} [s_{\hat{y}_2}(\mathbf{x}) - s_i(\mathbf{x})] < \delta \right) = 1,$$

where $\hat{y} := \operatorname{argmax}_{i \in [K]} s_i(\mathbf{x})$ and $\hat{y}_2 := \operatorname{argmax}_{i \neq \hat{y}, i \in [K]} s_i(\mathbf{x})$ denote the indices of the largest and second largest cosine similarities for an OOD input \mathbf{x} .

Theorem A.1. Given a pre-trained CLIP-like model $(\mathcal{T}, \mathcal{I})$ and a task with label set $\mathcal{Y}_{\text{in}} = \{y_1, y_2, \dots, y_K\}$. If $Q_{\mathbf{x}}$ satisfy Assumption A.1, Then there exists a constant $T = \frac{\lambda(K-1)(\lambda^{\text{wo}} + \delta - s_{\hat{y}_2})}{K\lambda - 1}$ such that for any temperature $\tau > T$, we have:

$$\text{FPR}(\tau, \lambda) \leq \text{FPR}^{\text{wo}}(\lambda^{\text{wo}}),$$

where $\text{FPR}(\tau, \lambda)$ is the false positive rate based on softmax scaling with temperature τ and threshold λ ; $\text{FPR}^{\text{wo}}(\lambda^{\text{wo}})$ is the false positive rate without softmax scaling based on threshold λ^{wo} . This suggests that applying softmax scaling with temperature results in superior OOD detection performance compared to without softmax scaling.

Proof. By definition, we express the false positive rate $\text{FPR}(\tau, \lambda)$ as follows,

$$\begin{aligned} \text{FPR}(\tau, \lambda) &= \mathbb{P}(G(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) = 1 \mid z = 0) \\ &= Q_{\mathbf{x}'}(G(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) = 1) \\ &= Q_{\mathbf{x}'}(p_{\hat{y}}(\mathbf{x}'; \tau) > \lambda) \\ &= Q_{\mathbf{x}'} \left(\frac{e^{s_{\hat{y}}(\mathbf{x}')/\tau}}{\sum_{j=1}^K e^{s_j(\mathbf{x}')/\tau}} > \lambda \right) \\ &= Q_{\mathbf{x}'} \left(\frac{1}{\lambda} > \sum_{i=1}^K \exp \left(\frac{s_i(\mathbf{x}') - s_{\hat{y}}(\mathbf{x}')}{\tau} \right) \right) \end{aligned}$$

²In practice, we observe that $s_i \in [0.1, 0.3]$ for CLIP with high probability.

By inequality $e^x \geq 1 + x$, we have,

$$Q_{\mathbf{x}'} \left(\frac{1}{\lambda} > \sum_{i=1}^K \exp \left(\frac{s_i(\mathbf{x}') - s_{\hat{y}}(\mathbf{x}')}{\tau} \right) \right) \leq Q_{\mathbf{x}'} \left(\frac{1}{\lambda} > \sum_{i=1}^K \left[1 + \frac{s_i(\mathbf{x}') - s_{\hat{y}}(\mathbf{x}')}{\tau} \right] \right)$$

This indicates

$$\begin{aligned} Q_{\mathbf{x}'} \left(\frac{1}{\lambda} > \sum_{i=1}^K \exp \left(\frac{s_i(\mathbf{x}') - s_{\hat{y}}(\mathbf{x}')}{\tau} \right) \right) &\leq Q_{\mathbf{x}'} \left(\frac{1}{\lambda} > \sum_{i=1}^K \left[1 + \frac{s_i(\mathbf{x}') - s_{\hat{y}}(\mathbf{x}')}{\tau} \right] \right) \\ &= Q_{\mathbf{x}'} \left(\sum_{i=1}^K (s_{\hat{y}}(\mathbf{x}') - s_i(\mathbf{x}')) > \left(K - \frac{1}{\lambda} \right) \tau \right) \end{aligned}$$

Since

$$\begin{aligned} \sum_{i=1}^K (s_{\hat{y}}(\mathbf{x}') - s_i(\mathbf{x}')) &= \sum_{i \neq \hat{y}} (s_{\hat{y}}(\mathbf{x}') - s_{\hat{y}_2}(\mathbf{x}') + s_{\hat{y}_2}(\mathbf{x}') - s_i(\mathbf{x}')) \\ &= \sum_{i \neq \hat{y}} (s_{\hat{y}}(\mathbf{x}') - s_{\hat{y}_2}(\mathbf{x}')) + \sum_{i \neq \hat{y}} (s_{\hat{y}_2}(\mathbf{x}') - s_i(\mathbf{x}')) \\ &= (K-1)(s_{\hat{y}}(\mathbf{x}') - s_{\hat{y}_2}(\mathbf{x}')) + \sum_{i \neq \hat{y}} (s_{\hat{y}_2}(\mathbf{x}') - s_i(\mathbf{x}')) \end{aligned}$$

By Assumption 3.1, we have

$$Q_{\mathbf{x}'} \left(\sum_{i=1}^K (s_{\hat{y}}(\mathbf{x}') - s_i(\mathbf{x}')) < (K-1)(s_{\hat{y}}(\mathbf{x}') - s_{\hat{y}_2}(\mathbf{x}')) + (K-1)\delta \right) = 1.$$

Therefore,

$$\begin{aligned} Q_{\mathbf{x}'} \left(\sum_{i=1}^K (s_{\hat{y}}(\mathbf{x}') - s_i(\mathbf{x}')) > \left(K - \frac{1}{\lambda} \right) \tau \right) &\leq Q_{\mathbf{x}'} \left(s_{\hat{y}}(\mathbf{x}') - s_{\hat{y}_2}(\mathbf{x}') > -\delta_2 + \left(K - \frac{1}{\lambda} \right) \frac{\tau}{K-1} \right) \\ &= Q_{\mathbf{x}'} (s_{\hat{y}}(\mathbf{x}') - s_{\hat{y}_2}(\mathbf{x}') > -\delta_2 + \lambda') \\ &= Q_{\mathbf{x}'} (s_{\hat{y}}(\mathbf{x}') > s_{\hat{y}_2}(\mathbf{x}') - \delta_2 + \lambda'), \end{aligned}$$

where $\lambda' = \left(K - \frac{1}{\lambda} \right) \frac{\tau}{K-1}$ is a monotonic function of λ (i.e., minimizing false positive rate *w.r.t.* λ is equivalent to minimizing *w.r.t.* λ' .)

For $\tau > 0$, we can rewrite the MCM score as

$$\begin{aligned} S_{\text{MCM}}(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) &= \max_{i \in [K]} \frac{e^{s_i(\mathbf{x}')/\tau}}{\sum_{j=1}^K e^{s_j(\mathbf{x}')/\tau}} = \frac{e^{s_{\hat{y}}(\mathbf{x}')/\tau}}{\sum_{j=1}^K e^{s_j(\mathbf{x}')/\tau}} \\ &= \frac{1}{1 + \sum_{j=1, j \neq \hat{y}}^K e^{(s_j(\mathbf{x}') - s_{\hat{y}}(\mathbf{x}'))/\tau}} \end{aligned}$$

As $\hat{y} := \operatorname{argmax}_{i \in [K]} s_i(\mathbf{x})$, $s_j(\mathbf{x}') - s_{\hat{y}}(\mathbf{x}') \leq 0$, $S_{\text{MCM}}(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I})$ is a monotonically decreasing function of τ , we have:

$$S_{\text{MCM}}(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) > \lim_{\tau \rightarrow \infty} \frac{1}{1 + \sum_{j=1, j \neq \hat{y}}^K e^{(s_j(\mathbf{x}') - s_{\hat{y}}(\mathbf{x}'))/\tau}} = \frac{1}{K}$$

Therefore by the definition of λ , we have $\lambda > \frac{1}{K}$, $\lambda' = \left(K - \frac{1}{\lambda} \right) \frac{\tau}{K-1} > 0$

For moderately large $\tau > T$ where $T = \frac{\lambda(K-1)(\lambda^{\mathbf{w}_0} + \delta - s_{\hat{y}_2})}{K\lambda - 1}$, we always have $s_{\hat{y}_2}(\mathbf{x}') - \delta + \lambda' > \lambda^{\mathbf{w}_0}$. Therefore, we obtain the following inequality,

$$\text{FPR}(\tau, \lambda) \leq Q_{\mathbf{x}'} (s_{\hat{y}}(\mathbf{x}') > s_{\hat{y}_2}(\mathbf{x}') - \delta_2 + \lambda') \leq Q_{\mathbf{x}'} (s_{\hat{y}}(\mathbf{x}') > \lambda^{\mathbf{w}_0}) := \text{FPR}^{\mathbf{w}_0}(\lambda^{\mathbf{w}_0}),$$

which means that the FPR without softmax scaling is larger than that with softmax scaling and a moderately large temperature. We show in Section 5 that the bound is indeed satisfied in practice with a large-scale ID dataset. \square

B Experimental Details

B.1 Software and Hardware

All methods are implemented in Pytorch 1.10. We run all OOD detection experiments on NVIDIA GeForce RTX-2080Ti GPU and use NVIDIA A100 GPU for fine-tuning CLIP and ViT.

B.2 Hyperparameters

The only hyperparameter in MCM is the (temperature) scaling factor τ . We use $\tau = 1$ by default unless otherwise specified. Our experiments suggest that MCM is insensitive to the scaling factor, where τ in a wide range of $[0.5, 100]$ shares similar performance.

B.3 Datasets

ImageNet-10 We create ImageNet-10 that mimics the class distribution of CIFAR-10 but with high-resolution images. It contains the following categories (with class ID): warplane (n04552348), sports car (n04285008), brambling bird, (n01530575), Siamese cat (n02123597), antelope (n02422699), Swiss mountain dog (n02107574), bull frog (n01641577), garbage truck (n03417042), horse (n02389026), container ship (n03095699).

ImageNet-20 For hard OOD evaluation with realistic datasets, we curate ImageNet-20, which consists of 20 classes semantically similar to ImageNet-10 (*e.g.*, dog (ID) vs. wolf (OOD)). The categories are selected based on the distance in the WordNet synsets [18]. Specifically, it contains the following categories: sailboat (n04147183), canoe (n02951358), balloon (n02782093), tank (n04389033), missile (n03773504), bullet train (n02917067), starfish (n02317335), spotted salamander (n01632458), common newt (n01630670), zebra (n01631663), frilled lizard (n02391049), green lizard (n01693334), African crocodile (n01697457), Arctic fox (n02120079), timber wolf (n02114367), brown bear (n02132136), moped (n03785016), steam locomotive (n04310018), space shuttle (n04266014), snowmobile (n04252077).

We hope the above two datasets will help future research on large-scale hard OOD detection. We provide a script for generating the datasets at <https://github.com/deeplearning-wisc/MCM>.

ImageNet-100 We randomly sample 100 classes from ImageNet-1k to curate ImageNet-100. To facilitate reproducibility, the script for generating the dataset and the class list are provided at <https://github.com/deeplearning-wisc/MCM>.

Conventional (non-spurious) OOD datasets Huang *et al.* [32] curate a diverse collection of subsets from iNaturalist [76], SUN [86], Places [96], and Texture [10] as large-scale OOD datasets for ImageNet-1k, where the classes of the test sets do not overlap with ImageNet-1k. We provide a brief introduction to each dataset as follows.

iNaturalist contains images in the natural world [76]. It has 13 super-categories and 5,089 sub-categories covering plants, insects, birds, mammals, and so on. We use the subset that contains 110 plant classes not overlapping with ImageNet-1k.

SUN stands for the Scene UNderstanding Dataset [86]. SUN contains 899 categories that cover more than indoor, urban, and natural places with or without human beings appearing. We use the subset which contains 50 natural objects not showing in ImageNet-1k.

Places is a large scene photographs dataset [96]. It contains photos that are labeled with scene semantic categories from three macro-classes: Indoor, Nature, and Urban. The subset we use is sampled from 50 categories that are not present in ImageNet-1k.

Texture stands for the Describable Textures Dataset [10]. It contains images of textures and abstracted patterns. As no categories overlap with ImageNet-1k, we use the entire dataset as in [32].

B.4 Baselines and sources of model checkpoints

For the Mahalanobis score [42], we use the feature embeddings without l_2 normalization as Gaussian distributions naturally do not fit hyperspherical features. Alternatively, one can normalize the embeddings first and then apply the Mahalanobis score.

For Fort *et al.* [19] in Table 2, we fine-tune the whole ViT model on the ID dataset. Specifically, we use the publicly available checkpoints from Hugging Face where the ViT model is pre-trained on ImageNet-21k and fine-tuned on ImageNet-1k. For example, the checkpoint for ViT-B is available at <https://huggingface.co/google/vit-base-patch16-224>.

For CLIP models, our reported results are based on checkpoints provided by Hugging Face for CLIP-B <https://huggingface.co/openai/clip-vit-base-patch16> and CLIP-L <https://huggingface.co/openai/clip-vit-large-patch14>. Similar results can be obtained with checkpoints in the codebase by OpenAI <https://github.com/openai/CLIP>. Note that for CLIP (RN50x4), which is not available in Hugging Face, we use the checkpoint provided by OpenAI.

C Spurious OOD Datasets

In general, spurious attributes refer to statistically informative features that co-exist with the majority of ID samples but do not necessarily capture cues related to the labels such as color, texture, background, etc [2, 3, 21, 87, 99]. A recent work [50] investigated a new type of hard OOD samples (called spurious OOD) that contain spurious or environmental features, but no object features related to the ID classes. A concrete example is shown in Figure 6, where images of birds co-occur frequently with either the land background or water background. Modern neural networks can spuriously rely on the image background (*e.g.*, water or land) for classification instead of learning to recognize the actual object [62]. Ming *et al.* [50] show that spurious OOD samples remain challenging for most common OOD detection methods based on pure vision models such as ResNet [23].

For ID dataset, we use Waterbirds [64], which combines bird photographs from CUB-200 [80] with water or land background images from PLACES [96]. For the spurious OOD dataset, we use the one created in [50] consisting of land and water background from Places [96].

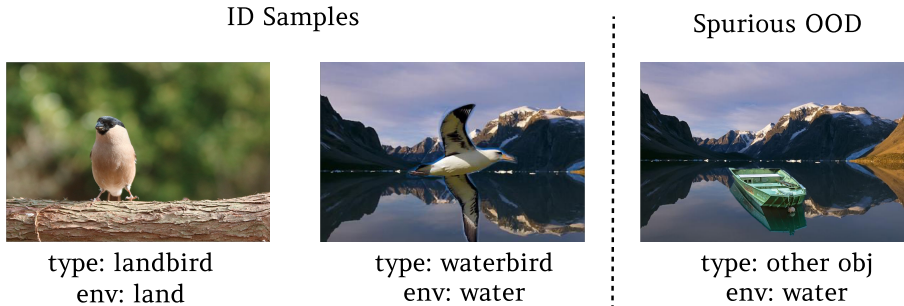


Figure 6: Illustration of spurious OOD samples for Waterbirds [64]. Images are taken from [50].

D ID Classification Accuracy

Table 7 shows the multi-class classification accuracy on ImageNet-1k for methods in Table 2.

Table 7: ID classification accuracy on ImageNet-1k (%)

Method	ID ACC
zero-shot	
MCM (CLIP-B/16)	67.01
MCM (CLIP-L/14)	73.28
w. fine-tuning	
MSP (CLIP-B/16)	79.39
MSP (CLIP-L/14)	84.12
Energy [48] (CLIP-B/16)	79.39
Energy [48] (CLIP-L/14)	84.12
Fort et al. [19] (ViT-B/16)	81.25
Fort et al. [19] (ViT-L/14)	84.05
MOS [32] (BiT)	75.16

E Implementation of CLIP-Based Baselines

E.1 Overview of Baselines

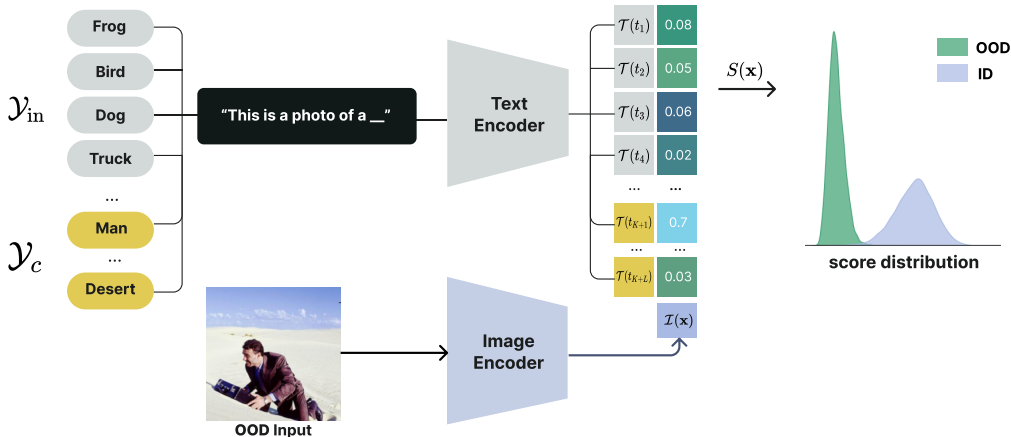


Figure 7: Zero-shot OOD detection with candidate OOD labels. The ID classification task is defined by a set of class labels \mathcal{Y}_{in} . With an additional set of candidate labels \mathcal{Y}_c that describes the contents of the input image, the OOD detection scoring function can be calculated by normalizing over the expanded space of cosine similarities.

We review two previous works on CLIP-based OOD detection [16, 19] in Figure 7, which derive the scoring function based on candidate OOD labels. For a given task with ID label set $\mathcal{Y}_{in} = \{y_1, y_2, \dots, y_K\}$ and candidate labels $\mathcal{Y}_c = \{y_{K+1}, y_{K+2}, \dots, y_{K+L}\}$, where ideally $\mathcal{Y}_{in} \cap \mathcal{Y}_c = \emptyset$, they construct a collection of text embeddings $\mathcal{T}(t_i), i \in \{1, 2, \dots, K + L\}$. Here, t_i is the text prompt “this is a photo of a $\langle y_i \rangle$ ” for a label y_i . For any test input image \mathbf{x} , we can calculate the label-wise matching score based on the cosine similarity between the image and text features: $s_i(\mathbf{x}) = \frac{\mathcal{I}(\mathbf{x}) \cdot \mathcal{T}(t_i)}{\|\mathcal{I}(\mathbf{x}')\| \cdot \|\mathcal{T}(t_i)\|}$. Therefore, a detection score can be derived as:

$$S(\mathbf{x}; \mathcal{Y}_{in}, \mathcal{Y}_c, \mathcal{T}, \mathcal{I}) = \sum_{i=1}^K \frac{e^{s_i(\mathbf{x})/\tau}}{\sum_{j=1}^{K+L} e^{s_j(\mathbf{x})/\tau}},$$

where $\tau > 0$ is the temperature scaling hyperparameter.

E.2 Obtaining OOD Candidate Labels

For the baseline methods, obtaining OOD candidate labels is a major challenge and limitation. Recently, [16] propose ZO-CLIP, where a transformer (decoder) based on the image encoder of CLIP is used to generate candidate labels. The transformer is trained from scratch on the COCO dataset [47] with simple teacher forcing algorithms. Although the decoder trained on COCO may work well on CIFAR (ID), it does not scale up to large-scale datasets such as ImageNet [11] where categories are not covered in COCO. As a result, [16] only test on small-scale datasets with common classes such as CIFAR (ID).

We improve the baseline by using a high-quality caption generator pre-trained on much larger datasets, which not only saves computational overhead but can potentially improve the quality of generated labels. The pipeline involves three components (see Figure 8):

- A caption generator. Given an input image, it generates a caption serving as the textual description of the input. In this work, we consider ClipCap [51], which uses GPT-2 [60] to generate captions based on CLIP’s image encoder. ClipCap is pre-trained on a much larger dataset Conceptual Captions [55] compared to COCO, which can be viewed as an enhanced version of the ZO-CLIP baseline [16]. The checkpoints are publicly available³.
- A syntactic parser. Given a caption, we extract noun objects using a parsing toolkit released by spaCy⁴. Those nouns can be used as candidate labels \mathcal{Y}_c of the input image.

³https://github.com/rmokady/CLIP_prefix_caption

⁴<https://spacy.io/models/en>

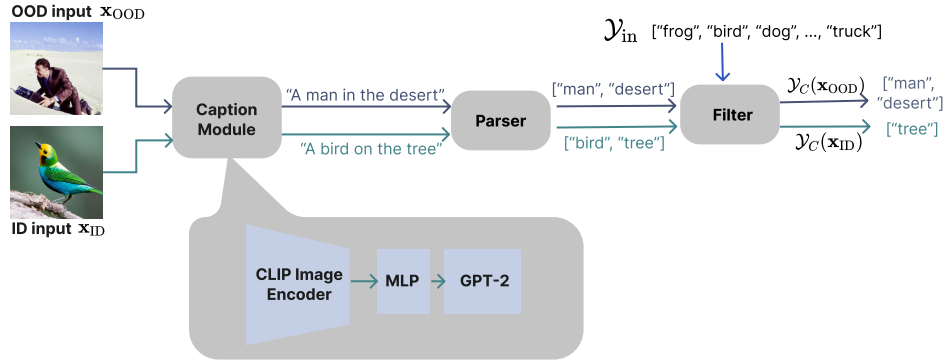


Figure 8: Improved pipeline to generate candidate OOD labels. It consists of three main components: a caption generator, a syntactic parser, and a filtering module to remove candidate labels that overlap with the ID label set.

- A filter module. Unlike [16], we further enhance the baseline by adopting a filtering technique to remove overlapping categories in \mathcal{Y}_C with ID labels \mathcal{Y}_{in} , which we detail below.

E.3 Label Filtering

Example. To illustrate the effects of filtering, we begin with a concrete example where ID labels are [“frog”, “bird”... “truck”], as shown in Figure 8. The generated labels (without filtering) of an ID input of a bird sitting on a tree are [“bird”, “tree”]. Therefore, $\mathcal{Y}_{in} \cup \mathcal{Y}_C = [“frog”, “bird” \dots “truck”, “bird”, “tree”]$. Ideally, the softmax probability distribution over the concatenated labels would be $[0, 0.5, 0, \dots, 0.5, 0]$ and by definition $S(\mathbf{x}) \approx 0.5$. However, if we filter the generated labels to eliminate nouns with similar meanings as ID, our concatenated labels would be [“frog”, “bird”... “truck”, “tree”] and the probability vector would be $[0, 1, 0, \dots, 0]$, which leads to a much higher score $S(\mathbf{x}) = 1$. In contrast, the generated labels for an OOD input with a caption “man in the desert” would be [“man”, “desert”]. The resulting probability vector would be $[0, 0, 0, \dots, 1, 0]$ and the score $S(\mathbf{x}) = 0$. Therefore, filtering makes it easier to separate ID inputs from OOD inputs (*c.f.* Figure 9).

String-based filtering. To implement the idea of filtering, we need a measurement of the similarity between the generated labels and ID labels. The simplest way is string-based filtering where a generated label is filtered if it matches any ID labels (in the string format), as in the case above.

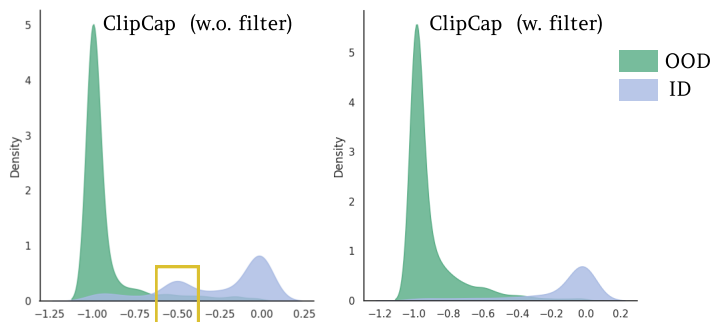


Figure 9: Score distributions for ImageNet-10 (ID) and iNaturalist (OOD) inputs. Simple string-based filtering alleviates the overlap between OOD inputs and ID inputs especially with scores around 0.5 (yellow rectangle), resulting in better ID-OOD separability.

F Alternative Scoring Functions

We explore the effectiveness of several alternative scoring functions:

- Entropy: the (negative) entropy of softmax scaled cosine similarities denoted as S_{entropy} ;
- Var: the variance of the cosine similarities denoted as S_{var} ;
- Scaled: the scaled difference between the largest and second-largest cosine similarities $S_{\text{scaled}} := e^{s_{\hat{y}_1(\mathbf{x})} - s_{\hat{y}_2(\mathbf{x})}}$ where $\hat{y}_1 := \operatorname{argmax}_{i \in [K]} s_i(\mathbf{x})$ and $\hat{y}_2 := \operatorname{argmax}_{i \neq \hat{y}_1, i \in [K]} s_i(\mathbf{x})$.

As shown in Table 8, MCM still gives the most promising results compared to the other three alternative scores across most OOD test sets.

Table 8: Comparison with other scaling functions (applied to inner products) on the large-scale benchmark ImageNet-1k (ID). We use CLIP-B/16 as the backbone.

Method	OOD Dataset									
	iNaturalist		SUN		Places		Texture		Average	
	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑
Entropy	84.44	63.50	93.79	62.54	94.10	64.15	97.16	58.98	92.37	62.29
Var	87.42	63.87	68.71	81.02	76.28	75.38	80.04	71.90	78.11	73.04
Scaled	89.06	72.26	89.06	70.81	89.08	69.66	89.56	68.17	89.19	70.22
MCM	30.91	94.61	37.59	92.57	44.69	89.77	57.77	86.11	42.74	90.77