

A UNIFIED APPROACH TOWARDS ACTIVE LEARNING AND OUT-OF-DISTRIBUTION DETECTION

Anonymous authors

Paper under double-blind review

ABSTRACT

When applying deep learning models in real-world scenarios, active learning (AL) strategies are crucial for identifying label candidates from a nearly infinite amount of unlabeled data. In this context, robust out-of-distribution (OOD) detection mechanisms are essential for handling data outside the target distribution of the application. However, current works investigate both problems separately. In this work, we introduce SISOM as the first unified solution for both AL and OOD detection. By leveraging feature space distance metrics SISOM combines the strengths of the currently independent tasks to solve both effectively. We conducted extensive experiments showing the problems arising when migrating between both tasks. In these evaluations SISOM underlined its effectiveness by achieving first place in two of the widely used OpenOOD benchmarks and second place in the remaining one. In AL, SISOM¹ outperforms others and delivers top-1 performance in three benchmarks.

1 INTRODUCTION

Large-scale deep learning models encounter several data-centric challenges during training and operation, particularly in real-world problems such as mobile robotic perception. On the one hand, these models require vast amounts of data and labels for training, driven by the uncontrolled nature of real-world tasks. On the other hand, even when trained with extensive data, these models can behave unpredictably when encountering samples that deviate significantly from the training data, known as out-of-distribution (OOD) data.

Active learning (AL) addresses the first limitation by guiding the selection of label candidates. In the traditional pool-based AL scenario (Settles, 2010), models start with a small labeled training set and can iteratively query data and its labels from an unlabeled data pool. The selection is based on model metrics such as uncertainty, diversity, or latent space encoding. One AL cycle concludes with the model being trained on the labeled subset, including the newly added samples.

The second challenge, dealing with unknown data during operation, is typically addressed by OOD detection. OOD detection distinguishes between in-distribution (InD) data used for training the model and OOD samples, which differ from the training distribution. Literature differentiates between near-OOD and far-OOD, which can be categorized by the type of distribution shifts occurring. Yang et al. (2022b) assumes near-OOD as a pure covariate shift while far-OOD often contains a semantic shift.

Over the whole life cycle of mobile robotic applications, which consists of training and operation phases, both challenges occur. Fig. 1 illustrates such a life cycle with both tasks. Given an amount of collected data, AL is applied for a label-efficient training, while OOD detection is employed to control

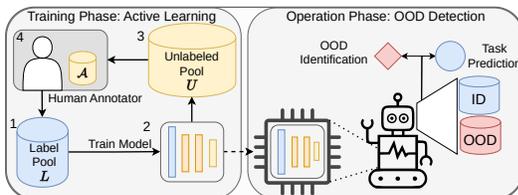


Figure 1: Real-world application life cycle comprising active learning in the training phase (left) and out-of-distribution detection in the operation phase (right).

¹SISOM will be published upon acceptance - for review <https://tinyurl.com/sisom-iclr>

the operation state, which is necessary for real-world operation domains. Existing works address these challenges separately, which can lead to diverging goals of AL and OOD methods. Additionally, addressing these tasks separately introduces significant overhead, especially for deployment and development like hyperparameter optimization or the training of auxiliary models.

From a method perspective similarities between AL and OOD detection are even more evident, specifically both methodologies utilize common metrics, such as uncertainty, latent space distances, and energy. In addition, a sample detected by such metrics can be, on the one hand, a novel AL sample that is insufficiently represented by the current training distribution. On the other hand, the sample can pose a covariate shift in an OOD setting. Considering both cases as depicted in Fig. 2 show an ambiguity and overlap of both sample categories. This raises the question if an examination of the ambiguity and relation between the respective samples can provide valuable insights for designing approaches for both tasks.

In our work, we examine the connection between both tasks and design a novel approach by leveraging mutual strengths providing an effective solution for both tasks. Specifically, we employ enriched feature space distances based on neural coverage to propose **Simultaneous Informative Sampling and Outlier Mining (SISOM)**, which create a symbiosis between AL and OOD detection. By exploiting the ambiguity of both tasks, SISOM effectively archives *top-1* performance in most OOD benchmarks and, at the same time, surpasses existing AL methods with *top-1* performance. With its joint approach, SISOM provides an efficient simplification for application life cycles by *eliminating an additional OOD detection design phase* and avoiding conflicting design goals. Additionally, SISOM provides a *novel latent space analysis for post-training latent space refinement* and a first-of-its-kind *self-balancing of uncertainty and diversity metrics*.

In summary, *our contributions* are as follows:

- We explore the entanglement of AL and OOD detection.
- We propose **Simultaneous Informative Sampling and Outlier Mining**, a novel method designed for both OOD detection *and* AL.
- We introduce a latent space analysis enabling an *optimization loop* for further *post-training latent space refinement* and a *self-balanced uncertainty diversity fusion*.
- In extensive experiments, we demonstrate SISOM effectiveness in AL *and* OOD benchmarks.

2 PRELIMINARIES

Active Learning: AL is a subfield of machine learning designed to reduce the number of required labels by querying a set of new samples \mathbb{A} of a query size q in a cyclic process. Let \mathcal{X} represent a set of samples and \mathcal{Y} a set of labels. AL starts with an initially labeled pool \mathbb{L} , containing data samples with features \mathbf{x} and corresponding label y , and an unlabeled pool \mathbb{U} where only \mathbf{x} is known. However, y can be queried from a human oracle. We further assume that \mathbb{L} and \mathbb{U} are samples from a distribution Ω . In each cycle, a model f is trained such that $f : \mathcal{X}_L \rightarrow \mathcal{Y}_L$. This model then selects new samples from \mathbb{U} based on a query strategy $Q(\mathbf{x}, f)$, which utilizes (intermediate) model outputs. As a result, the newly annotated set \mathbb{A} is added to the labeled pool \mathbb{L}^{i+1} and removed from the unlabeled pool \mathbb{U}^{i+1} .

Out of Distribution Detection: Ancillary, OOD detection assumes a model $f : \mathcal{X}_L \rightarrow \mathcal{Y}_L$ trained on our training data $\{\mathbf{x}, y\} \in \mathbb{L}$ which have been sampled from the distribution Ω . During evaluation or inference, a model f encounters data samples $\tilde{\mathbf{x}}$ from a distribution Θ and Ω , where $\Omega \cap \Theta = \emptyset$ and $\tilde{\mathbf{x}} \notin \mathbb{L}$. Data sampled from Ω are referred to as InD data, while samples from Θ are referred to as OOD data. Based on the trained model f , a metric S is used to determine whether a sample x is

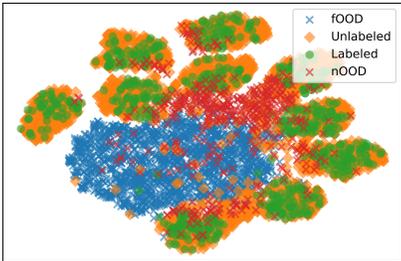


Figure 2: TSNE plot of unlabeled and OOD data compared to labeled data for CIFAR-10 as InD with 20% labeled, tiny ImageNet as near-OOD and SVHN as far-OOD.

108 sampled from Ω or Θ .

$$109 \quad G(\mathbf{x}, f) = \begin{cases} \text{InD} & \text{if } S(\mathbf{x}; f) \geq \lambda \\ \text{OOD} & \text{if } S(\mathbf{x}; f) < \lambda \end{cases} \quad (1)$$

111 OOD detection is further categorized into near- and far-OOD (Zhang et al., 2023). Far-OOD refers to
112 completely unrelated data, such as comparing MNIST (LeCun et al., 1998) to CIFAR-100 (Krizhevsky
113 et al., 2009), while CIFAR-10 (Krizhevsky et al., 2009) to CIFAR-100 would be considered as near-
114 OOD. OpenOOD (Yang et al., 2022b) ranks near-OOD detection as more challenging.

116 3 RELATED WORK

117 Given the disentanglement of fields, we review the related work individually.

118 **Active Learning:** AL mainly considers the pool-based and stream-based scenario (Settles, 2010),
119 where data is either queried from a pool in a data center or a stream on the fly. For deep learning, the
120 majority of current research deals with pool-based AL (Ren et al., 2021). However, further scenarios
121 have been evaluated by Schmidt & Günnemann (2023) and Schmidt et al. (2024). Independent of the
122 scenarios, samples are selected either by prediction uncertainty, latent space diversity, or auxiliary
123 models. A majority of the uncertainty-based methods rely on sampling - like Monte Carlo Dropout
124 (Gal & Ghahramani, 2016) - or employ ensembles (Beluch et al., 2018; Lakshminarayanan et al.,
125 2017). To additionally ensure batch diversity Kirsch et al. (2019) used the joint mutual information.
126 The uncertainty concepts have been employed and further developed for major computer vision tasks,
127 including object detection (Feng et al., 2019; Schmidt et al., 2020), 3D object detection (Hekimoglu
128 et al., 2022; Park et al., 2023), and semantic segmentation (Huang et al., 2018). One of the few works
129 breaking the gap between both tasks (Shukla et al., 2022) modified an OOD detection method for pose
130 estimation. Mukhoti et al. (2023) proposed an uncertainty baseline based on spectral convolutions and
131 Gaussian mixture models, which shows effectiveness on AL and OOD detection compared to other
132 uncertainty approaches. In contrast, diversity-based approaches aim to select key samples to cover
133 the whole dataset. Sener & Savarese (2018) proposed to choose a CoreSet of the latent space using a
134 greedy optimization. Yehuda et al. (2022) selected samples having high coverage in a fixed radius for
135 low data regimes. Mishal & Weinshall (2024) extends the approach for more data regimes dynamic
136 strategy mixing. Ash et al. (2020) enriched the latent space dimensions to the dimensions of the
137 gradients and included uncertainty in this way. The concept of combining uncertainty with diversity
138 has been further refined for 3D object detection (Yang et al., 2022a; Luo et al., 2023). Liang et al.
139 (2022) combined different diversity metrics for the same task. In semantic segmentation, Surprise
140 Adequacy (Kim et al., 2020) has been employed to measure how surprising a model finds a new
141 instance. Besides the metric-based approach, the selection can also be made by auxiliary models
142 mimicking diversity and uncertainty. These approaches range from loss estimation (Yoo & Kweon,
143 2019), autoencoder-based approaches (Sinha et al., 2019; Zhang et al., 2020; Kim et al., 2021) and
144 graph models (Caramalau et al., 2021), to teacher-student approaches (Peng et al., 2021; Hekimoglu
145 et al., 2024).

146 **Out-of-Distribution Detection:** To facilitate a fair comparison and evaluation of OOD methods,
147 benchmarking frameworks like OpenOOD (Yang et al., 2022b; Zhang et al., 2023) have been intro-
148 duced, which categorizes the methods into preprocessing methods altering the training process and
149 postprocessing methods being applied after training. Preprocessing techniques include augmenting
150 training data like mixing (Zhang et al., 2018; Tokozume et al., 2018) different samples or applying
151 fractals to images (Hendrycks et al., 2022). Postprocessing approaches include techniques of ma-
152 nipulations on neurons and weights of the trained network, such as filtering for important neurons
153 (Ahn et al., 2023; Djuricic et al., 2022), or weights (Sun & Li, 2022), or clipping neuron values to
154 reduce OOD-induced noise (Sun et al., 2021). Logit-based approaches encompass the model output
155 to estimate uncertainties using temperature-scaling (Liang et al., 2018), modified entropy scores
156 (X. Liu, 2023), energy scores (Liu et al., 2020; Elfein et al., 2021) or ensembles (Arpit et al., 2022).
157 Other methods rely on distances in the feature space, such as the Mahalanobis distance between InD
158 and OOD samples (Lee et al., 2018), consider the gradients after a forwardpass (Liang et al., 2018;
159 Hsu et al., 2020; Huang et al., 2021; Schwinn et al., 2021), estimate densities (Charpentier et al.,
160 2020; 2022) or k nearest neighbor on latent space distances (Sun et al., 2022). A different branch
161 operates on the features directly and evaluates properties like the Norm (Yu et al., 2023) or performs
rank reductions via SVD (Song et al., 2022). NAC (Liu et al., 2024) combined gradient information
with a density approach, where a probability density function over InD samples is estimated.

OpenSet Active Learning: The emerging field of OpenSet AL considers both tasks in one cycle, assuming the AL pool is polluted by OOD samples. Existing approaches (Ning et al., 2022; Park et al., 2022; Yang et al., 2023; Safaei et al., 2024) address both tasks with *separate* modules containing auxiliary models. None of the works investigates the correlation of AL and OOD samples. As both tasks are considered decouples with uncorrelated modules, this field is orthogonal to our examination of correlation and entanglement. We believe that this field profits from the joint consideration of AL and OOD samples as well as an examination of their ambiguity.

While various works exist in OOD and AL, both tasks are considered independent. Even in OpenSet AL, the tasks are considered by independent method components. Some uncertainty methods are evaluated on both tasks but limit their evaluation to the uncertainty domain. Current state-of-the-art approaches are often specified for one task. In addition, the application life cycle consideration is unexplored.

4 METHODOLOGY

To address both AL and OOD detection tasks in a unified method to simplify real-world applications, we need to first understand the goals of these two tasks. AL aims to identify and select samples that are beneficial for training and increase the models performance. These samples typically position themselves between the existing clusters in the latent space or near the decision boundaries. OOD detection targets the identification of data outside the training data and, therefore, outside the known clusters in latent space. Given the definition of far- and near-OOD, near-OOD is closer to InD data and located close to the decision boundaries and in between the existing clusters. Liu & Qin (2024) recently showed that OOD is generally closer to the decision boundary than InD confirming this hypothesis. Fig. 2 depicts this consideration showing the overlap of interesting unlabeled data and (near-)OOD data.

To target these overlapping regions we design a method focusing on the latent space regions between the clusters. To do so SISOM employs an enlarged feature space Coverage (1) and increases expressiveness by weighting important neurons in a Feature Enhancement (2). Based on this feature representation, we refine the AL selection and the InD and OOD border by using an inner-to-outer class Distance Ratio (3), guiding it to unexplored and decision boundary regions. As feature space distances are prone to poorly defined latent space representations, we introduce Feature Space Analysis (4) providing a self-deciding fusion of our distance metric with an uncertainty-based energy score. Optionally, our previous analysis enables us to optimize the Sigmoid Stepness (5), providing a further refinement of the feature space representations from (2). An overview is depicted in Fig. 3.

(1) Coverage: We aim to identify the regions of the samples that are interesting and unexplored for AL as well as OOD samples in latent space. To do so, we rely on an informative latent space covering as much information as possible.

To increase the information gain we cover the full network and define the feature space representation of an input sample \mathbf{x} as a concatenation of the latent space of multiple layers h_j in a set of selected layers H in Eq. (2). This approach follows the procedures of neural coverage (Kim et al., 2019; Liu et al., 2024) and is contrasting to most diversity-based AL approaches (Sener & Savarese, 2018; Ash et al., 2020), which use a single layer.

$$\mathbf{z} = h_1(\mathbf{x}) \oplus \dots \oplus h_j(\mathbf{x}) \oplus \dots \oplus h_n(\mathbf{x}) \quad (2)$$

Given the feature space \mathbf{z} , we further denote \mathbb{Z}_U as a set of feature space representations of unlabeled samples from \mathbb{U} , while \mathbb{Z}_L denotes the set of representations of all labeled samples \mathbb{L} .

(2) Feature Enhancement: To enhance the expressiveness of our defined latent space we introduce a weighting of individual layers. Prior research (Huang et al., 2021; Liu et al., 2024) have demonstrated that the gradients of neurons with respect to the KL divergence of the model’s output and a uniform distribution encapsulate valuable information for OOD detection.

We apply the technique to improve the features further and enrich these by representing the individual contribution of each neuron i , denoted as g_i . This gradient describes each neuron’s contribution to the actual output being different from the uniform distribution. A low value suggests that the neuron has little influence on the prediction of a given input sample. Conversely, if the value is high, the respective neuron is crucial for the decision process.

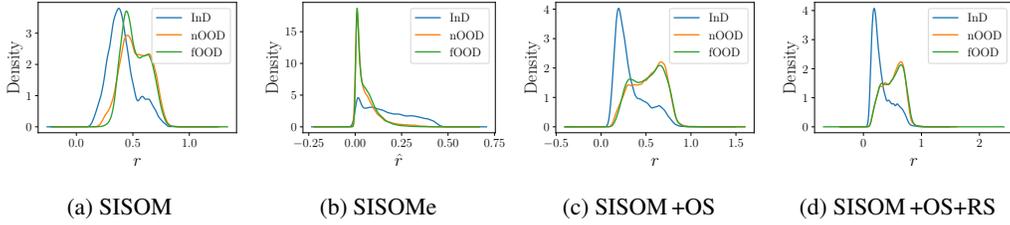


Figure 4: Density plots for SISOM with energy, Optimal Sigmoid Steepness (OS) and Reduced Subset Selection (RS) on CIFAR-100 with near-OOD (nOOD) and far-OOD (fOOD) as defined in OpenOOD.

Thus, the gradient vector can be interpreted as a saliency weighting for the activation values in the feature space to support separability. In detail, we compute the gradient of the Kullback-Leibler (KL) divergence between an uniform distribution u and the softmax output distribution $f(\mathbf{x})$ for an input \mathbf{x} :

$$\mathbf{g}_i = \frac{\partial D_{KL}(u||f(\mathbf{x}))}{\partial \mathbf{z}_i}. \quad (3)$$

We incorporate the calculated saliency to create a weighted feature representation forming the enhanced feature space with the sigmoid function σ :

$$\tilde{\mathbf{z}} = \sigma(\mathbf{z} \odot \mathbf{g}). \quad (4)$$

The resulting gradient-weighted feature representation effectively prioritizes the most influential neurons for each input. This facilitates the identification of inputs activating atypical influence patterns, which is significant for AL as well as OOD detection. A qualitative analysis demonstrating the effect of the feature enrichment is given in Appendix A.3.

(3) Distance Ratio: After we defined and enhanced our latent space we design our metric to identify the respective samples. Contrasting to other works in the latent space domain for AL and OOD detection (Sener & Savarese, 2018) which rely on simple distance metrics, we take inspiration from complex distance metrics (Kim et al., 2019) for detecting adversarial examples.

We assume the location of important samples in between the existing clusters in latent space. While samples closer to these clusters, like near-OOD or AL samples close to the decision boundary, are more important, far-OOD samples and exotic AL samples should not be omitted. To identify samples in these regions, we rely on a distance quotient between inner-class and outer-class distances.

The inner-class distance d_{in} is defined as the minimal feature space distance to a known sample of the same class c as the predicted pseudo-class of the given sample. The outer-class distance d_{out} represents the minimal feature space distance to a known sample of a different class than the sample’s pseudo-class.

$$d_{in} = \min_{\mathbf{z}' \in \mathbb{Z}_L(c'=c)} \|\tilde{\mathbf{z}} - \tilde{\mathbf{z}}'\|_2 \quad (5) \quad d_{out} = \min_{\mathbf{z}' \in \mathbb{Z}_L(c' \neq c)} \|\tilde{\mathbf{z}} - \tilde{\mathbf{z}}'\|_2 \quad (6)$$

The distance is computed on the gradient-enhanced feature space $\tilde{\mathbf{z}}$ defined in Eq. (4) with \mathbf{z}' describing the nearest sample from the set of known samples \mathbb{Z}_L .

In many state-of-the-art works on AL, computationally expensive distance calculations are often present (Sener & Savarese, 2018; Ash et al., 2020; Caramalau et al., 2021). To make our approach

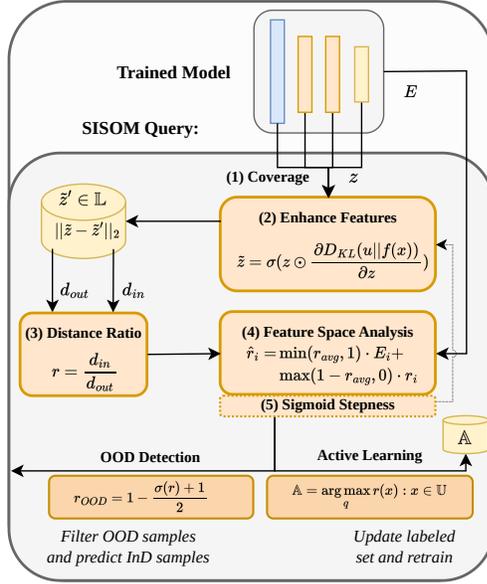


Figure 3: SISOM framework for OOD detection and AL combined.

more efficient for AL and feasible for large-scale OOD detection tasks, we select a representative subset $\mathbb{T} \subset \mathbb{Z}_L$ as a comparison set, thereby significantly reducing computational overhead. We modify the Probcover (Yehuda et al., 2022) approach to select class-wise samples, maximizing coverage within a sphere with a fixed radius in the feature space. The effect of this subset selection is further investigated in Section 5.3.

Our SISOM score r reflects the distance between each neuron’s weighted feature representation in the latent space and the nearest sample of the predicted class relative to the closest distance to a sample from a different class:

$$r = \frac{d_{in}}{d_{out}}. \quad (7)$$

An extended comparison of the different distance metrics and their ability to separate InD and OOD is shown in Appendix A.3, while a SISOM is depicted in Fig. 4a.

For AL we select the q samples with the highest distance ratio r , with q being the AL query size:

$$\mathbb{A} = \operatorname{argmax}_q r(\mathbf{x}) : \mathbf{x} \in \mathbb{U}. \quad (8)$$

For OOD Detection, we map the distance ratios r to an interval $[0; 1]$ with the strictly monotonically decreasing function:

$$r_{OOD} = 1 - \frac{\sigma(r) + 1}{2}. \quad (9)$$

(4) Feature Space Analysis: Having a well-defined latent space is crucial for SISOM to attain optimal performance. Furthermore, we hypothesize that techniques relying on feature space metrics are more dependent on feature space separation than uncertainty-based methods. This dependency is important for SISOM as it utilizes a quotient of feature space metrics. Nevertheless, obtaining a well-defined and separable latent space may pose challenges in specific contexts and tasks.

To estimate the separability of feature space, we compute the average distance ratio r_{avg} using Eq. (4) and Eq. (7) for the known set as:

$$r_{avg} = \frac{1}{|\mathbb{L}|} \sum_{\tilde{\mathbf{z}} \in \mathbb{L}} \frac{d_{in}(\tilde{\mathbf{z}})}{d_{out}(\tilde{\mathbf{z}})} = \frac{1}{|\mathbb{L}|} \sum_{\mathbf{z} \in \mathbb{L}} \frac{d_{in}(\sigma(\mathbf{z} \odot \mathbf{g}))}{d_{out}(\sigma(\mathbf{z} \odot \mathbf{g}))}. \quad (10)$$

A lower r_{avg} value indicates better separation of the samples in the enhanced feature space, implying that samples of the same class are relatively closer together than samples of different classes. To mitigate possible performance disparities of SISOM in difficult separable domains, we introduce a novel self-deciding process for the sampling method, which utilizes the feature separation score r_{avg} as follows:

$$\hat{r}_i = \min(r_{avg}, 1) \cdot E_i + \max(1 - r_{avg}, 0) \cdot r_i. \quad (11)$$

The so created \hat{r} combines our SISOM score from Eq. (7) with the uncertainty-based energy score $E(\mathbf{x}) = -\log \sum_{i=1}^c \exp(f(\mathbf{x})_i)$ based on the model’s output logits $f(\mathbf{x})$.

Depending on whether $r_{avg} \rightarrow 1$ or $r_{avg} \rightarrow 0$, the created score \hat{r}_i relies more on either the energy score or the distance ratio r_i . If $r_{avg} \rightarrow 1$, indicating poorly separated classes, \hat{r}_i relies more on the energy score. Conversely, if $r_{avg} \rightarrow 0$, suggesting a well-separated feature space, \hat{r}_i relies more on the distance ratio. A density outline of our combined approach SISOMe is given in Fig. 4b. Alternatively, one can replace r_{avg} with a tuneable hyperparameter in Eq. (11).

(5) Sigmoid Steepness: Since Eq. (10) depends on the sigmoid function defined in Eq. (4), the sigmoid function has a large influence on the enhanced feature space $\tilde{\mathbf{z}}$. An additional hyperparameter α can influence the sigmoid function’s steepness. As \mathbf{z} is concatenated from different layers in Eq. (2), the sigmoid can be applied to each layer j individually. This allows for a more nuanced control over the influence of each neuron’s contribution to the final decision and so influences the separability of the feature space. We define the sigmoid using the steepness parameter α as:

$$\sigma_j(\mathbf{x}) = \frac{1}{1 + e^{-\alpha_j \mathbf{x}}}; \quad \{\alpha_j : h_j \in \mathbf{z} \ \forall j\}. \quad (12)$$

Relating to Eq. (4), the set α of steepness parameters of the sigmoid function for each layer h_j , determines the degree of continuity or discreteness of the features within that layer. By applying a

layerwise sigmoid, Eq. (4) is formulated as follows:

$$\begin{aligned} \tilde{\mathbf{z}} &= \sigma_1(h_1(\mathbf{x}) \odot \mathbf{g}_{i,1}) \oplus \cdots \oplus \sigma_j(h_j(\mathbf{x}) \odot \mathbf{g}_{i,j}) \oplus \\ &\quad \cdots \oplus \sigma_n(h_n(\mathbf{x}) \odot \mathbf{g}_{i,n}), \end{aligned} \quad (13)$$

with $\mathbf{g}_{i,j} = \frac{\partial D_{KL}(u||f(\mathbf{x}))}{\partial h_{j,i}}; \quad \forall j.$

Following this consideration we can select α values which optimize the feature space separability metric r_{avg} from Eq. (10) by minimizing $\alpha_{opt} = \arg \min_{\alpha} r_{avg}(\alpha)$. Besides the quantitative assessment of our Feature Space Analysis and Sigmoid Steepness in Section 5.3, the influence of the Sigmoid Steepness is shown in Fig. 4c.

5 EXPERIMENTS

To evaluate our proposed method, we conducted a comprehensive assessment of SISOM on both tasks AL and OOD detection individually. We consider compound tasks like Openset AL as out of scope as existing works address the sub-task by individual components, while SISOM showcases the ambiguity of both task sample characteristics. The experiments’ details, settings, and results are presented in Section 5.1 and Section 5.2, respectively. We further conduct an ablation study in Section 5.3. We utilized the standard pool-based AL scenario (Settles, 2010) for AL. For OOD detection, we followed the widely used OpenOOD benchmarking framework (Yang et al., 2022b; Zhang et al., 2023).

In the AL experiments, we compared our method against several baselines, including **CoreSet** (Sener & Savarese, 2018), **CoreGCN** (Caramalau et al., 2021), **Random Badge** (Ash et al., 2020), and **Loss Learning** (Yoo & Kweon, 2019). Additionally, we adapted the **NAC** (Liu et al., 2024) method from OOD detection to AL to assess the transferability from OOD to AL.

For OOD detection experiments, we employed the implementation provided by the OpenOOD framework when available. We also followed the experimental setup and datasets for near- and far-OOO detection. The baselines used for validation include **NAC** (Liu et al., 2024), **Ash** (Djurisic et al., 2022), **KNN** (Sun et al., 2022), **Odin** (Hsu et al., 2020), **ReAct** (Sun et al., 2021), **MSP** (Hendrycks & Gimpel, 2016), **Energy** (Liu et al., 2020), **Dice** (Sun & Li, 2022), **RankFeat** (Yu et al., 2023), **FeatureNorm** (Song et al., 2022) and **GEN** (X. Liu, 2023). Moreover, we tested the **CoreSet** (Sener & Savarese, 2018) AL method to verify the transferability from AL to OOD. Our focus was on methods that use the cross-entropy training scheme to maintain a fair comparison and ensure compatibility post-AL.

5.1 ACTIVE LEARNING

We followed the most common AL benchmark settings and datasets, including the CIFAR-10 (Krizhevsky et al., 2009), CIFAR-100 (Krizhevsky et al., 2009), and SVHN (Netzer et al., 2011) datasets paired with a ResNet18 (He et al., 2016) model. We assessed the network’s performance by measuring accuracy relative to the amount of data used. The plots include markers to indicate the selection steps. As suggested by (Yoo & Kweon, 2019; Ash et al., 2020), we start with an initial pool size of 1,000 labeled samples for CIFAR-10 and SVHN. In each AL cycle, the model can query 1,000 additional samples from an unlabeled pool, which are then labeled and added to the labeled pool for the subsequent cycle. Due to the larger number of classes in CIFAR-100, we increased the selection size to 5,000. Detailed parameters and settings are available in Appendix B.1.

In the CIFAR-10 benchmark depicted in Fig. 5a, SISOM exhibits swift progress and maintains consistent performance from the outset. It consistently outperforms other methods, achieving the highest performance differential in all selection cycles and is only eclipsed by SISOMe especially in early cycles. Furthermore, as the sample size increases, our method maintains its superiority over Learning Loss and CoreSet. NAC does not demonstrate superior performance compared to Random.

After examining SISOM in datasets with a limited number of classes, we examine the AL setup on the larger CIFAR-100 dataset and report the results in Fig. 5b. In this setting, all methods are less stable in its ranking compared to the other dataset, reflecting the increased difficulty of the dataset. The complexity of the dataset requires more data for the model to perform effectively. While in the

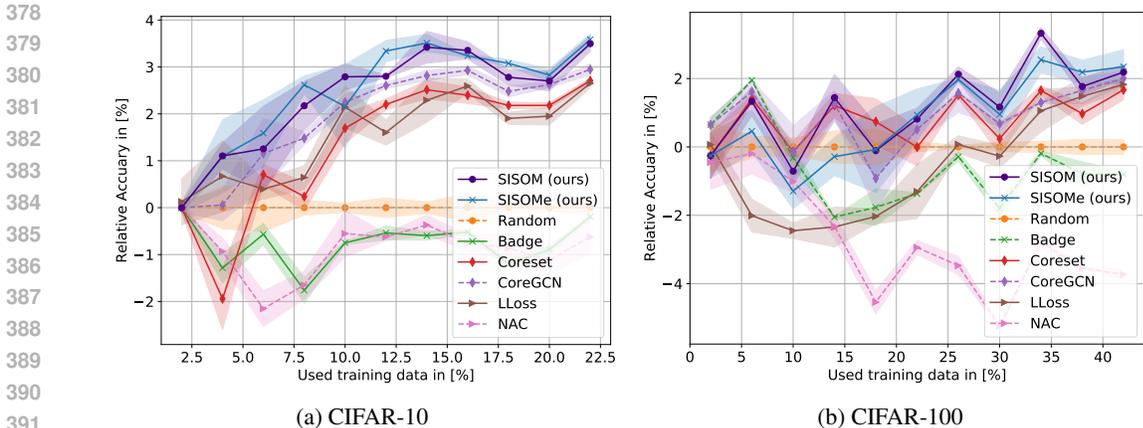


Figure 5: Comparison of different active learning methods on CIFAR-10, SVHN and CIFAR-100 with indicated standard errors.

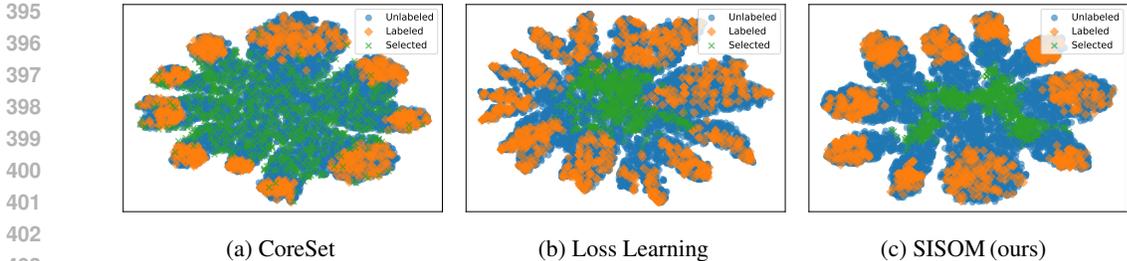


Figure 6: T-SNE feature space comparison of Loss Learning, CoreSet and SISOM for SVHN on cycle 1. SISOM effectively targets the areas in-between the clusters.

early stages, pure diversity-based methods are in the lead, SISOM gains velocity in the last selection steps and achieves the highest performance difference only in the last step SISOMe is more effective.

Following the experiments on CIFAR-10 and CIFAR-100 we conducts experiments on SVHN and report them in Appendix A.1.

In conclusion of the AL experiments, SISOM reached state-of-the-art performance and surpasses other methods across all three datasets, demonstrating its viability for AL. While in the early stages, SISOM falls behind other approaches for CIFAR-100, in following selection cycles with more training data it outperforms them. We hypothesize that the early cycles had a poorly separated feature space, causing this issue.

5.2 OUT-OF-DISTRIBUTION DETECTION

Following our evaluation of SISOM on classic AL benchmarks, we utilize the OpenOOD framework to evaluate its performance on the OOD detection task. We stick to the recommended benchmarks on CIFAR-10 (Krizhevsky et al., 2009), CIFAR-100 (Krizhevsky et al., 2009), and ImageNet 1k (Deng et al., 2009), and we provide evaluation values for both near- and far-OOD detection. The assignment of datasets to near and far categories follows the framework’s suggestions and is reported with additional settings in Appendix B.2. In addition, we benchmarked the life cycle setting in Appendix A.2. The framework ranks methods based on their AUROC performance and provides checkpoints for fair post-processor validation.

Firstly, we examine the performance on the CIFAR-10 benchmark and show the results in Table 1a. SISOMe and SISOM achieve the highest AUROC score for near-OOD data, respectively. SISOMe surpasses SISOM in all metrics. For far-OOD, SISOM ranks third after NAC, while SISOMe secures the first place. This is noteworthy as NAC underperformed in the AL task, even when compared to methods suffering from batch diversification, which underlines the non-triviality of migrating between both tasks out of the box.

Table 1: OOD benchmark for CIFAR-10, CIFAR-100 and ImageNet1k with Cross-Entropy training setting and dataset according to OpenOOD sorted by Near-OOD performance.

	(a) CIFAR-10			(b) CIFAR-100			(c) ImageNet 1k				
Post-processor	OOD AUROC		ID Acc.	Post-processor	OOD AUROC		ID Acc.	Post-processor	OOD AUROC		ID Acc.
	Near-OOD	Far-OOD			Near-OOD	Far-OOD			Near-OOD	Far-OOD	
SISOMe	91.76	94.74	95.06	Gen	81.31	79.68	77.25	SISOMe	78.59	89.04	76.18
SISOM	<u>91.40</u>	94.50	95.06	SISOMe	<u>80.96</u>	79.8	77.25	ASH	<u>78.17</u>	95.74	76.18
NAC	90.93	<u>94.60</u>	95.06	Energy	80.91	79.77	77.25	ReAct	77.38	93.67	76.18
KNN	90.64	92.96	95.06	ReAct	80.77	80.39	77.25	SISOM	77.33	88.01	76.18
CoreSet	90.34	92.85	95.06	MSP	80.27	77.76	77.25	GEN	76.85	89.76	76.18
GEN	88.20	91.35	95.06	KNN	80.18	82.4	77.25	KLM	76.64	87.6	76.18
MSP	88.03	90.73	95.06	ODIN	79.9	79.28	77.25	Energy	76.03	89.50	76.18
Energy	87.58	91.21	95.06	SISOM	79.42	77.91	77.25	MSP	76.02	85.23	76.18
ReAct	87.11	90.42	95.06	DICE	79.38	80.01	77.25	ODIN	74.75	89.47	76.18
FeatureNorm	85.52	95.59	95.06	ASH	78.2	<u>80.58</u>	77.25	DICE	73.07	90.95	76.18
ODIN	82.87	87.96	95.06	KLM	76.56	76.24	77.25	NAC	71.73	<u>94.66</u>	76.18
RankFeat	79.46	75.87	95.06	CoreSet	75.69	79.53	77.25	KNN	71.1	90.18	76.18
KLM	79.19	82.68	95.06	NAC	72.00	86.56	77.25	FeatureNorm	67.57	91.13	76.18
DICE	78.34	84.23	95.06	RankFeat	61.88	67.10	77.25	RankFeat	50.99	53.93	76.18
ASH	75.27	78.49	95.06	FeatureNorm	47.87	80.99	77.25	Coreset	-	-	76.18

In the OpenOOD CIFAR-100 benchmark Table 1b the best far-OOD method shows the worst near-OOD performance, while for CIFAR-10, methods performed almost equally well on both near- and far-OOD. SISOMe ranks as the second-best method for near-OOD and repeatedly beats the individual metrics, SISOM and Energy. This is an interesting finding since, in contrast to CIFAR-10, energy achieves better performance than SISOM among the individual metrics on CIFAR-100. This supports our hypothesis that by considering the average ratio r_{avg} as a proxy for feature space separation, we obtain stronger performances in both well-separated and poorly-separated feature spaces.

The third benchmark suggested by OpenOOD is ImageNet 1k, which contains more classes and is a much larger dataset than the previous ones. In the results depicted in Table 1c, SISOMe and SISOM achieved first and fourth-best scores on near-OOD, with SISOMe showing strong performance for far-OOD. Interestingly, the NAC method, which was the second-best in CIFAR-10, ranks much lower, and KNN, the third-best method in CIFAR-10, ranks last. Meanwhile, ASH, which ranks first in this benchmark, is last in the CIFAR-10 benchmark.

To evaluate the life cycle perspective we conducted additional experiments using the AL models in A.2.

Overall benchmarks, SISOMe is the only approach, being consistently under the top three ranks, and even secured first place in two of them. Excluding SISOMe, SISOM achieved one top-three ranking and one top-one ranking. Notably, our method performs relatively better on near-OOD data than on far-OOD data. This is understandable, as the ratio between inner and outer class distance is higher for data close to the training data distribution, while the quotient is lower for far-OOD. Additionally, near-OOD is closer to the data of interest for AL selection. According to (Yang et al., 2022b), near-OOD is considered the more challenging task and is more likely to occur in real-world applications. Thus, higher performance on near-OOD may be preferred in practice.

5.3 ABLATIONS STUDIES

In an ablation study, we qualitatively examine the latent space assumptions for AL as well as the effect of unsupervised feature space analysis and reduce labeled set \mathbb{T} . A study of the individual components of SISOM is given in Appendix A.3.

AL Latent Space: To validate the assumptions made in Section 4, we examine the configuration of the latent space of our selection in the AL experiments. The objective of our method is to select samples in the decision boundary region for the AL case. In Fig. 6, we compare CoreSet and Loss Learning with SISOM. It can be observed that CoreSet, as intended, exhibits high diversity in unseparated regions. The pseudo-uncertainty-based Loss Learning method is more concentrated in its selection but fails to diversify the selection across all decision boundaries. In contrast, SISOM, as shown in Fig. 6c, focuses on the decision boundary while successfully covering the entire area between the unseparated samples. This demonstrates the effectiveness of our method in addressing the challenges of both AL and OOD detection.

Optimal Sigmoid Steepness: In our feature space analysis in Section 4, we derived r_{avg} in Eq. (10) as a proxy for the feature space separability. Due to the distance concept of SISOM, we hypothesize that it works better in well-separated feature spaces. To examine this, we conduct a random search for different α sets and record the different r_{avg} values. To reduce the search space, we follow the premise postulated in Section 4 that generally, deeper layers require a steeper sigmoid curve, i.e., a higher α_j value due to the nature of the features captured within these layers.

After computing every r_{avg} value for each combination of α , we select the α_{opt} set that minimizes r_{avg} . Formally, this can be written as:

$$\alpha_{\text{opt}} = \arg \min_{\alpha} r_{\text{avg}}(\alpha)$$

In Table 2, an optimized set α_{opt} is marked with OS. As it can be seen, a set with better feature space separation leads to increased performance for CIFAR-100 and ImageNet, partly confirming our hypothesis. In CIFAR-10 however, the original set of parameters yields the best results. One explanation might be that, in CIFAR-10, the different classes are already well separated, such that optimization on this separation yields no improvement and leads to an overfitting behavior.

Reduced Subset Selection: For larger datasets, distance-based approaches like CoreSet (Sener & Savarese, 2018) or (Ash et al., 2020) suffer from huge computational efforts, which is problematic for OOD detection, too. In Section 4, we suggested to use a reduced subset \mathbb{T} of the comparison set \mathbb{Z}_L , selecting class-wise samples with the most neighbors in a given radius. For each dataset, we select a total of 10% of the samples for each class, drastically increasing inference speed. We compare the effect of our reduced subset selection (RS) in Table 2 and highlight it qualitatively in Fig. 4d. A comparison of the preprocessing steps for SISOM in Table 2 indicates that the AUROC near-OOD score has improved for all datasets. It can be observed that preselection enhances feature space separability based on the r_{avg} column. This also strengthens our hypothesis from the previous subsection. For ImageNet and CIFAR-100, the combination of feature analysis and preselection results in the best performance, for CIFAR-10 the additional feature space analysis did not improve the performance. By taking the low r_{avg} into account, the chosen values could have reduced the space too much, leading to an overfitting behavior. All parameters are given in Appendix B.3.

Table 2: Ablation Study on Optimal Sigmoid Steepness (OS) and Reduced Subset Selection (RS) on Near OOD Benchmarks.

Method	ImageNet		CIFAR 100		CIFAR 10	
	AUROC _n	r_{avg}	AUROC _n	r_{avg}	AUROC _n	r_{avg}
SISOM	77.21	0.270	75.93	0.33	<u>91.33</u>	0.26
SISOM, OS	<u>77.4</u>	0.266	<u>79.56</u>	0.19	90.37	0.099
SISOM, RS	77.33	0.249	76.07	0.31	91.40	0.24
SISOM, OS, RS	77.37	0.245	79.69	0.18	90.54	0.086

6 CONCLUSION

We proposed SISOM, the first approach designed to solve OOD detection and AL jointly, providing an effective simplification in real-world application life cycles by eliminating an OOD design phase and avoiding conflicting goals of AL and OOD detection. By weighting latent space features with KL divergence of the neuron activations and relating them to the latent space clusters of the different classes SISOM achieves state-of-the-art performance in both tasks. In addition, SISOM provides a novel feature space analysis scheme enabling a post-training feature space refinement as well as a self-guided uncertainty and diversity fusion introduced as SISOMe. In the famous OpenOOD benchmarks SISOM archives the *top-1* performance in *two of the three benchmarks* and the second place in the remaining one. For active learning, SISOM surpasses state-of-the-art approaches in three different benchmarks. While current state-of-the-art approaches are highly specialized for either AL or OOD detection, SISOM solves both tasks with the same approach. Underlined by these results, SISOM effectively addresses real-world applications, like environment sensing, which usually suffers from label costs during training and high unlabeled data availability as well as out-of-distribution samples during inference.

In future work, we plan to combine the two tasks that are currently separated as independent steps. Enabling continuous AL during inference while filtering out-of-distribution data can significantly enhance the model’s performance after the initial selection phase.

540 REPRODUCEABILITY STATEMENT

541

542 To ensure reproducibility, we conducted all experiments with the same fixed seeds, which are reported
 543 in the training procedure in the appendix. We used the exact parameter setting of the OpenOOD
 544 benchmark for the OOD experiments. Moreover, the code is released in the benchmark form with
 545 available configurations.

546

547 ETHICS STATEMENT

548

549 With our research, we address the challenges of real-world and mobile (robotic) applications. While
 550 the common usage of robots or real-world applications does not pose ethical concerns, these fields
 551 pose the risk of misuse.

552

553 REFERENCES

554

555 Yong Hyun Ahn, Gyeong-Moon Park, and Seong Tae Kim. LINE: Out-of-Distribution Detection by
 556 Leveraging Important Neurons. In *Proceedings of the IEEE/CVF Conference on Computer Vision
 557 and Pattern Recognition (CVPR)*, 2023.

558

559 Devansh Arpit, Huan Wang, Yingbo Zhou, and Caiming Xiong. Ensemble of averages: Improving
 560 model selection and boosting performance in domain generalization. In *Proceedings of the
 561 International Conference on Neural Information Processing Systems (NeurIPS)*, 2022.

562

563 Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep
 564 batch active learning by diverse, uncertain gradient lower bounds. In *Proceedings of the Interna-
 565 tional Conference on Learning Representations (ICLR)*, 2020.

566

567 William H Beluch, Tim Genewein, Andreas Nürnberger, and Jan M Köhler. The power of ensembles
 568 for active learning in image classification. In *Proceedings of the IEEE/CVF Conference on
 569 Computer Vision and Pattern Recognition (CVPR)*, 2018.

570

571 Julian Bitterwolf, Maximilian Mueller, and Matthias Hein. In or out? fixing imagenet out-of-
 572 distribution detection evaluation. In *Proceedings of the International Conference on Machine
 573 Learning (ICML)*, 2023.

574

575 Razvan Caramalau, Binod Bhattarai, and Tae-Kyun Kim. Sequential graph convolutional network
 576 for active learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
 577 Recognition (CVPR)*, 2021.

578

579 Bertrand Charpentier, Daniel Zügner, and Stephan Günnemann. Posterior network: Uncertainty esti-
 580 mation without ood samples via density-based pseudo-counts. In *Proceedings of the International
 581 Conference on Neural Information Processing Systems (NeurIPS)*, 2020.

582

583 Bertrand Charpentier, Oliver Borchert, Daniel Zügner, Simon Geisler, and Stephan Günnemann.
 584 Natural posterior network: Deep bayesian uncertainty for exponential family distributions. In
 585 *International Conference on Machine Learning (ICML)*, 2022.

586

587 M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and and A. Vedaldi. Describing Textures in the Wild.
 588 In *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014.

589

590 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale
 591 hierarchical image database. In *Proceedings of the IEEE/CVF Conference on Computer Vision
 592 and Pattern Recognition (CVPR)*, 2009.

593

594 Andrija Djurisic, Nebojsa Bozanic, Arjun Ashok, and Rosanne Liu. Extremely simple activation
 595 shaping for out-of-distribution detection. In *Proceedings of the International Conference on
 596 Learning Representations (ICLR)*, 2022.

597

598 Sven Elflein, Bertrand Charpentier, Daniel Zügner, and Stephan Günnemann. On out-of-distribution
 599 detection with energy-based models. *arXiv*, 2107.08785, 2021.

- 594 Di Feng, Xiao Wei, Lars Rosenbaum, Atsuto Maki, and Klaus Dietmayer. Deep active learning for
595 efficient training of a lidar 3d object detector. In *Proceedings of the IEEE Intelligent Vehicles
596 Symposium (IV)*, 2019.
- 597 Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model
598 uncertainty in deep learning. In *Proceedings of the International Conference on Machine Learning
599 (ICML)*, 2016.
- 600 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
601 recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
602 Recognition (CVPR)*, 2016.
- 603 Aral Hekimoglu, Michael Schmidt, Alvaro Marcos-Ramiro, and Gerhard Rigoll. Efficient active
604 learning strategies for monocular 3d object detection. In *Proceedings of the IEEE Intelligent
605 Vehicles Symposium (IV)*, 2022.
- 606 Aral Hekimoglu, Michael Schmidt, and Alvaro Marcos-Ramiro. Monocular 3d object detection with
607 lidar guided semi supervised active learning. In *Proceedings of the IEEE/CVF Winter Conference
608 on Applications of Computer Vision (WACV)*, 2024.
- 609 Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution
610 examples in neural networks. In *Proceedings of the International Conference on Learning
611 Representations (ICLR)*, 2016.
- 612 Dan Hendrycks, Andy Zou, Mantas Mazeika, Leonard Tang, Bo Li, Dawn Song, and Jacob Steinhardt.
613 Pixmix: Dreamlike pictures comprehensively improve safety measures. In *Proceedings of the
614 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- 615 Yen-Chang Hsu, Yilin Shen, Hongxia Jin, and Zsolt Kira. Generalized odin: Detecting out-of-
616 distribution image without learning from out-of-distribution data. In *Proceedings of the IEEE/CVF
617 Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- 618 Po Yu Huang, Wan Ting Hsu, Chun Yueh Chiu, Ting Fan Wu, and Min Sun. Efficient uncertainty
619 estimation for semantic segmentation in videos. In *Proceedings of the European Conference on
620 Computer Vision (ECCV)*, 2018.
- 621 Rui Huang, Andrew Geng, and Yixuan Li. On the importance of gradients for detecting distributional
622 shifts in the wild. In *Proceedings of the International Conference on Neural Information Processing
623 Systems (NeurIPS)*, 2021.
- 624 Jinhan Kim, Robert Feldt, and Shin Yoo. Guiding Deep Learning System Testing Using Surprise
625 Adequacy. In *Proceedings of the International Conference on Software Engineering (ICSE)*, 2019.
- 626 Jinhan Kim, Jeongil Ju, Robert Feldt, and Shin Yoo. Reducing DNN labelling cost using surprise
627 adequacy: an industrial case study for autonomous driving. In *Proceedings of the Joint European
628 Software Engineering Conference and Symposium on the Foundations of Software Engineering
629 (ESEC/FSE)*, 2020. ISBN 978-1-4503-7043-1.
- 630 Kwanyoung Kim, Dongwon Park, Kwang In Kim, and Se Young Chun. Task-aware variational
631 adversarial active learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
632 Pattern Recognition (CVPR)*, 2021.
- 633 Andreas Kirsch, Joost Van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch
634 acquisition for deep bayesian active learning. In *Proceedings of the International Conference on
635 Neural Information Processing Systems (NeurIPS)*, 2019.
- 636 Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Learning multiple layers of features from
637 tiny images. Technical report, Canadian Institute for Advanced Research, 2009. URL [http:
638 //www.cs.toronto.edu/~kriz/cifar.html](http://www.cs.toronto.edu/~kriz/cifar.html).
- 639 kuangliu. pytorch-cifar, 2021. URL <https://github.com/kuangliu/pytorch-cifar>.
640 GitHub repository.

- 648 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive
649 uncertainty estimation using deep ensembles. In *Proceedings of the International Conference on*
650 *Neural Information Processing Systems (NeurIPS)*, 2017.
- 651 Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. Technical Report 7, Stanford
652 Computer Vision Lab, 2015.
- 654 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to
655 document recognition. In *Proceedings of the IEEE*, 1998.
- 656 Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting
657 out-of-distribution samples and adversarial attacks. In *Proceedings of the International Conference*
658 *on Neural Information Processing Systems (NeurIPS)*, 2018.
- 660 Shiyu Liang, Yixuan Li, and R Srikant. Enhancing the reliability of out-of-distribution image detection
661 in neural networks. In *Proceedings of the International Conference on Learning Representations*
662 *(ICLR)*, 2018.
- 663 Zhihao Liang, Xun Xu, Shengheng Deng, Lile Cai, Tao Jiang, and Kui Jia. Exploring diversity-based
664 active learning for 3d object detection in autonomous driving. *arXiv*, 2205.07708, 2022.
- 666 Litian Liu and Yao Qin. Fast decision boundary based out-of-distribution detector. In *ICML*, 2024.
- 667 Weitang Liu, Xiaoyun Wang, John D. Owens, and Yixuan Li. Energy-based Out-of-distribution
668 Detection. In *Proceedings of the International Conference on Neural Information Processing*
669 *Systems (NeurIPS)*, October 2020.
- 671 Yibing Liu, Chris Xing Tian, Haoliang Li, Lei Ma, and Shiqi Wang. Neuron Activation Cover-
672 age: Rethinking Out-of-distribution Detection and Generalization. In *The Twelfth International*
673 *Conference on Learning Representations (ICLR)*, 2024.
- 674 Alejandro López-Cifuentes, Marcos Escudero-Vinolo, Jesús Bescós, and Álvaro García-Martín.
675 Semantic-aware scene recognition. *Pattern Recognition*, 102, 2020. Publisher: Elsevier.
- 677 Yadan Luo, Zhuoxiao Chen, Zijian Wang, Xin Yu, Zi Huang, and Mahsa Baktashmotlagh. Exploring
678 active 3d object detection from a generalization perspective. In *Proceedings of the International*
679 *Conference on Learning Representations (ICLR)*, 2023.
- 680 Inbal Mishal and Daphna Weinshall. Dcom: Active learning for all learners. *Arxiv*, 7 2024.
- 682 Jishnu Mukhoti, Andreas Kirsch, Joost Van Amersfoort, Philip H S Torr, and Yarin Gal. Deep
683 deterministic uncertainty: A new simple baseline. In *CVPR*, 2023. Introduction of a novel
684 uncertainty method and tested them on AL, OOD.
- 685 Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. Reading
686 digits in natural images with unsupervised feature learning. In *Proceedings of the International*
687 *Conference on Neural Information Processing Systems (NeurIPS) Workshop on Deep Learning*
688 *and Unsupervised Feature Learning*, 2011.
- 690 Kun-Peng Ning, Xun Zhao, Yu Li, and Sheng-Jun Huang. Active learning for open-set annotation. In
691 *CVPR*, 1 2022.
- 692 Dongmin Park, Yooju Shin, Jihwan Bang, Youngjun Lee, Hwanjun Song, and Jae-Gil Lee. Meta-
693 query-net: Resolving purity-informativeness dilemma in open-set active learning. In *Neurips*, 10
694 2022.
- 696 Younghyun Park, Wonjeong Choi, Soyeong Kim, Dong-Jun Han, and Jaekyun Moon. Active learning
697 for object detection with evidential deep learning and hierarchical uncertainty aggregation. In
698 *Proceedings of the International Conference on Learning Representations (ICLR)*, 2023.
- 699 Fengchao Peng, Chao Wang, Jianzhuang Liu, Zhen Yang Noah, and Ark Lab. Active learning for
700 lane detection: A knowledge distillation approach. In *Proceedings of the IEEE/CVF International*
701 *Conference on Computer Vision (ICCV)*, 2021.

- 702 Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Brij B Gupta, Xiaojiang Chen,
703 and Xin Wang. A survey of deep active learning. *ACM computing surveys (CSUR)*, 54(9):1–40,
704 2021.
- 705 Bardia Safaei, Vibashan VS, Celso M. de Melo, and Vishal M. Patel. Entropic open-set active
706 learning. In *AAAI*, 12 2024. URL <http://arxiv.org/abs/2312.14126>. Two different
707 entropies... ;br/;Closed Set Entropy is cal-;br/;culated based on the outputs of K class-aware binary
708 clas-;br/;sifiers (BC) trained on DL.;br/;Distance-based;br/;Entropy is utilized to prioritize the
709 selection of samples that;br/;stand apart from distributions of unknown classes;br/;;br/;Appling
710 Triplet Loss;br/;;br/;Related Work.;br/;OpenMax with an open class.
- 711 Sebastian Schmidt and Stephan Günnemann. Stream-based active learning by exploiting temporal
712 properties in perception with temporal predicted loss. In *Proceedings of the British Machine Vision
713 Conference (BMVC)*, 2023.
- 714 Sebastian Schmidt, Qing Rao, Julian Tatsch, and Alois Knoll. Advanced active learning strategies for
715 object detection. In *Proceedings of the IEEE Intelligent Vehicles Symposium (IV)*, 2020.
- 716 Sebastian Schmidt, Lukas Stappen, Leo Schwinn, and Stephan Günnemann. Generalized synchro-
717 nized active learning for multi-agent-based data selection on mobile robotic systems. *IEEE
718 Robotics and Automation Letters*, 9(10):8659–8666, 2024. doi: 10.1109/LRA.2024.3444670.
- 719 Leo Schwinn, An Nguyen, René Raab, Leon Bungert, Daniel Tenbrinck, Dario Zanca, Martin Burger,
720 and Bjoern Eskofier. Identifying untrustworthy predictions in neural networks by geometric
721 gradient analysis. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2021.
- 722 Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set
723 approach. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 8
724 2018.
- 725 Burr Settles. Active learning literature survey. Technical report, University of Wisconsin–Madison,
726 2010.
- 727 Megh Shukla, Roshan Roy, Pankaj Singh, Shuaib Ahmed, and Alexandre Alahi. VL4Pose: Active
728 Learning Through Out-Of-Distribution Detection For Pose Estimation. In *Proceedings of the
729 British Machine Vision Conference (BMVC)*, 2022.
- 730 Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell. Variational adversarial active learning. In
731 *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- 732 Yue Song, Nicu Sebe, and Wei Wang. Rankfeat: Rank-1 feature removal for out-of-distribution
733 detection. In *Proceedings of the International Conference on Neural Information Processing
734 Systems (NeurIPS)*, 2022. URL <https://github.com/KingJamesSong/RankFeat>.
- 735 Yiyou Sun and Yixuan Li. DICE: Leveraging Sparsification for Out-of-Distribution Detection. In
736 *Proceedings of the European Conference on Computer Vision (ECCV)*, November 2022.
- 737 Yiyou Sun, Chuan Guo, and Yixuan Li. React: Out-of-distribution detection with rectified activa-
738 tions. In *Proceedings of the International Conference on Neural Information Processing Systems
739 (NeurIPS)*, 2021.
- 740 Yiyou Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. Out-of-distribution detection with deep nearest
741 neighbors. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2022.
- 742 Yuji Tokozume, Yoshitaka Ushiku, and Tatsuya Harada. Between-class learning for image classifica-
743 tion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition
744 (CVPR)*, 2018.
- 745 Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam,
746 Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In
747 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
748 2018.

- 756 Sagar Vaze, Kai Han, Andrea Vedaldi, and Andrew Zisserman. Open-set recognition: A good
757 closed-set classifier is all you need. In *Proceedings of the International Conference on Learning*
758 *Representations (ICLR)*, 2021.
- 759 Haoqi Wang, Zhizhong Li, Litong Feng, and Wayne Zhang. Vim: Out-of-distribution with virtual-
760 logit matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
761 *Recognition (CVPR)*, 2022.
- 762 weiaicunzai. pytorch-cifar100, 2022. URL [https://github.com/weiaicunzai/
764 pytorch-cifar100](https://github.com/weiaicunzai/pytorch-cifar100). GitHub repository.
- 765 C. Zach X. Liu, Y. Lochman. GEN: Pushing the Limits of Softmax-Based Out-of-Distribution Detec-
766 tion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*
767 *(CVPR)*, 2023.
- 768 Chenhongyi Yang, Lichao Huang, and Elliot J. Crowley. Plug and play active learning for object
769 detection. *arXiv*, 2211.11612, 11 2022a.
- 770
771 Jingkang Yang, Pengyun Wang, Dejian Zou, Zitang Zhou, Kunyuan Ding, WenXuan Peng, Haoqi
772 Wang, Guangyao Chen, Bo Li, Yiyu Sun, Xuefeng Du, Kaiyang Zhou, Wayne Zhang, Dan
773 Hendrycks, Yixuan Li, and Ziwei Liu. OpenOOD: Benchmarking generalized out-of-distribution
774 detection. In *Proceedings of the International Conference on Neural Information Processing*
775 *Systems (NeurIPS) Datasets and Benchmarks Track*, 2022b.
- 776
777 Yang Yang, Yuxuan Zhang, Xin Song, and Yi Xu. Not all out-of-distribution data are harmful to
778 open-set active learning. In *Proceedings of the International Conference on Neural Information*
779 *Processing Systems (NeurIPS)*, 2023.
- 780 Ofer Yehuda, Avihu Dekel, Guy Hacohen, and Daphna Weinshall. Active learning through a covering
781 lens. In *Proceedings of the International Conference on Neural Information Processing Systems*
782 *(NeurIPS)*, 2022.
- 783
784 Donggeun Yoo and In So Kweon. Learning loss for active learning. In *Proceedings of the IEEE/CVF*
785 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- 786
787 Yeonguk Yu, Sungho Shin, Seongju Lee, Changhyun Jun, and Kyoobin Lee. Block selection method
788 for using feature norm in out-of-distribution detection. In *Proceedings of the IEEE/CVF Conference*
789 *on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- 790
791 Beichen Zhang, Liang Li, Shijie Yang, Shuhui Wang, Zheng-Jun Zha, and Qingming Huang. State-
792 relabeling adversarial active learning. In *Proceedings of the IEEE/CVF Conference on Computer*
793 *Vision and Pattern Recognition (CVPR)*, 2020.
- 794
795 Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical
796 risk minimization. In *Proceedings of the International Conference on Learning Representations*
797 *(ICLR)*, 2018.
- 798
799 Jingyang Zhang, Jingkang Yang, Pengyun Wang, Haoqi Wang, Yueqian Lin, Haoran Zhang, Yiyu
800 Sun, Xuefeng Du, Kaiyang Zhou, Wayne Zhang, Yixuan Li, Ziwei Liu, Yiran Chen, and Hai
801 Li. OpenOOD v1.5: Enhanced benchmark for out-of-distribution detection. In *Proceedings of*
802 *the International Conference on Neural Information Processing Systems (NeurIPS) Workshop on*
803 *Distribution Shifts: New Frontiers with Foundation Models*, 2023.
- 804
805
806
807
808
809

A ADDITIONAL EXPERIMENTS

In this section, we present additional experiments for SISOM.

A.1 ACTIVE LEARNING - SVHN

Following the CIFAR experiments settings, we depict the results for the SVHN experiments in Fig. 7. Similar to the CIFAR-10 results, our method maintained high performance, but the method differences shrink with the easier the dataset. In the last cycle, SISOM reaches the highest performance, with a margin over other methods. As for CIFAR-10, NAC did not perform well in the data selection. Given that SVHN’s 10 classes are numbers, it is easier than the more diverse CIFAR-10 benchmark dataset. This can be observed by an overall reduced performance gap between the methods compared to CIFAR-10.

A.2 OUT-OF-DISTRIBUTION LIFE CYCLE

To evaluate the effectiveness of SISOM in a life cycle setting, we utilized the models after the AL cycle for an OOD benchmark. In Table 3 we used the same setting as for the benchmark CIFAR-10 experiments with the similar near- and far-OOD. It should be noted that while openOOD is open to deploy different checkpoints, modifying the InD data access is more challenging and remains unchanged. In Table 3 SISOMe archived the top performance, making it suitable for the full application life cycle.

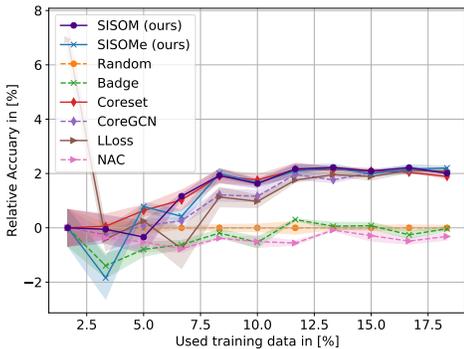


Table 3: OOD benchmark for CIFAR-10 using the AL checkpoints of SISOM.

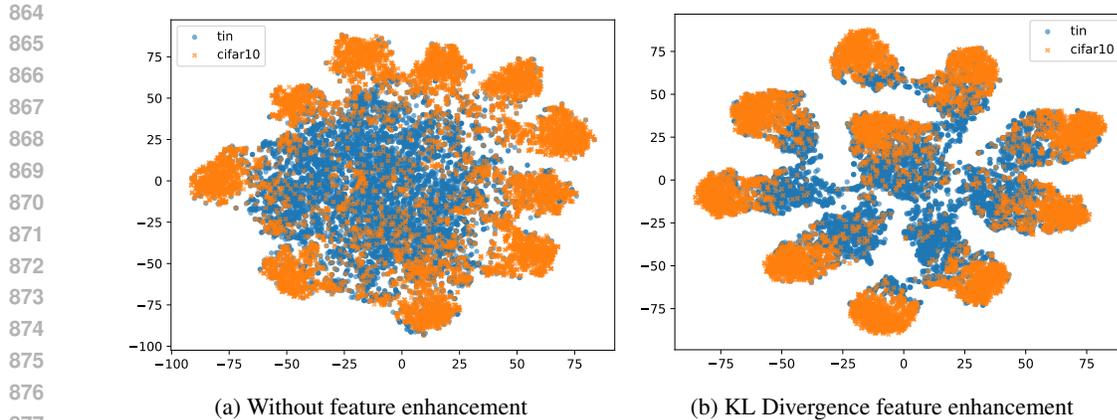
Postprocessor	OOD AUROC		ID Acc.
	Near-OOD	Far-OOD	
SISOMe	86.84	88.39	89.73
ReAct	86.84	87.72	89.73
GEN	85.43	86.04	89.73
MSP	84.37	84.85	89.73
ASH	83.39	87.33	89.61
NAC	82.26	85.06	89.73
RankFeat	60.20	56.73	60.84

Figure 7: Comparison of different active learning methods on SVHN with indicated standard errors.

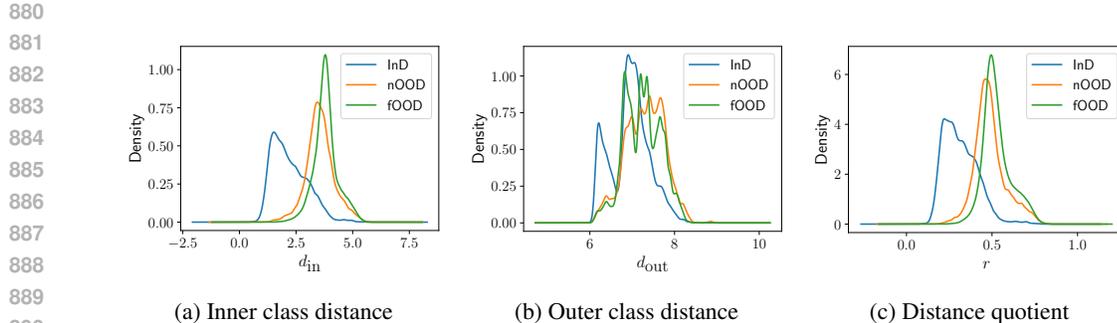
A.3 FEATURE SPACE ASSIGNMENTS

In this section, we highlight the influence of major components of our methods on the ability to separate InD and OOD data. In Fig. 8, we display the influence of the KL divergence gradient with a T-SNE analysis on CIFAR-10 (Krizhevsky et al., 2009) as InD and Tiny ImageNet (tin) (Le & Yang, 2015) as near-OOD. Without feature enhancement, the latent space is much harder to separate, and tin is distributed all over the latent space as shown in Fig. 8a. In contrast, the latent space with KL divergence enhances features, is much more separated, and has a clearer decision boundary to the near classes as indicated in Fig. 8b.

In addition to the previously presented density plots, we show the inner and outer distance together with the distance quotient of SISOM in Fig. 9 for CIFAR-10. Fig. 9a shows the inner class, indicating small inner class distances leading to a good separability for the InD data. On the other hand, the outer class distance in Fig. 9b provides a good separable peak for InD data, but a portion of InD overlaps with OOD data. The combined distance quotient shows the increased separability of the different InD and OOD sets as depicted in Fig. 9c.



878 Figure 8: T-SNE comparison of the latent space for OOD detection with and without KL-Divergence
879 feature enrichment.



891 Figure 9: Density plots for the inner class distance, outer class distance, and the distance quotient of
892 SISOM for CIFAR-10 with near-OOD (nOOD) and far-OOD (fOOD) as defined in OpenOOD..

895 B EXPERIMENTAL DETAILS

896
897 In this section, we provide experiment details to support the reproducibility of results by providing
898 the used parameters².

900 B.1 ACTIVE LEARNING EXPERIMENTS

901
902 In active learning experiments, we used a ResNet18 (He et al., 2016) model, with the suggested
903 modifications of (Yoo & Kweon, 2019) presented in a CIFAR benchmark repository (kuangliu, 2021),
904 which replaced the kernel of the first convolution with a 3×3 kernel. Additionally, we used an SGD
905 optimizer with a learning rate of 0.1 and multistep scheduling at 60, 120, and 160, decreasing the
906 learning rate by a factor of 10, which are reported benchmark parameters for CIFAR-100 (weiaicunzai,
907 2022). For SVHN and CIFAR-10 we used a learning rate of 0.025 and a cosine scheduler as suggested
908 by Yehuda et al. (2022). For the construction of the feature space, we used the layers after the 4
909 blocks of ResNet with the following sigmoid values:

- 910 • **CIFAR-10**
911 Adaptive Average Pooling Layer: 50,
912 Sequential Layer 3: 10,
913 Sequential Layer 2: 1,
914 Sequential Layer 1: 0.05.
- 915 • **CIFAR-100/SVHN**
916 Adaptive Average Pooling Layer: 1,
917

²Code will published upon acceptance, for review <https://tinyurl.com/sisom-iclr>

918 Sequential Layer 3: 0.1,
 919 Sequential Layer 2: 0.1,
 920 Sequential Layer 1: 0.1
 921

922 B.2 OUT-OF-DISTRIBUTION EXPERIMENTS

923
 924 In the OOD experiments, we report the mean of the three different seeds employed in the standard
 925 setting of the OpenOOD (Yang et al., 2022b) framework with ResNet18 for CIFAR-10 and CIFAR-
 926 100. For Imagenet, we use the sole ResNet50 torchvision checkpoint provided in the standard
 927 settings. We utilized the near- and far-OOD assignments suggested by the benchmark listed below.
 928 We followed the official tables of OpenOOD’s benchmark and reported the mean without the standard
 929 deviation. For the CIFAR-100 experiment, instead of using the automated r_{avg} value to balance
 930 between r and E from Eq. (11), we set $r_{\text{avg}} = 0.8$ for SISOMe based on a hyperparameter study. In
 931 the benchmark tables, we reported for SISOM the best values matching the best values of the ablation
 932 study modifications. Furthermore, we follow the suggested sigmoid values (Liu et al., 2024) for
 933 CIFAR-10 and ImageNet. For CIFAR-100, we choose values that minimize Eq. (10). A detailed
 934 overview of the sigmoid values for the 4 blocks of ResNet18 and ResNet50 for all experiments is
 935 provided below:

- 936 • **CIFAR-10**
 Adaptive Average Pooling Layer: 100,
 Sequential Layer 3: 1000,
 Sequential Layer 2: 0.001,
 Sequential Layer 1: 0.001
- 941 • **CIFAR-100**
 Adaptive Average Pooling Layer: 1,
 Sequential Layer 3: 0.1,
 Sequential Layer 2: 0.1,
 Sequential Layer 1: 0.1
- 946 • **ImageNet:**
 Adaptive Average Pooling Layer: 3000,
 Sequential Layer 3: 300,
 Sequential Layer 2: 0.01,
 Sequential Layer 1: 1

950
 951 OOD dataset assignment:

- 952 • **CIFAR-10**
 Near-OOD: CIFAR-100 (Krizhevsky et al., 2009), Tiny ImageNet (Le & Yang, 2015)
 Far-OOD: MNIST (LeCun et al., 1998), SVHN (Netzer et al., 2011), Textures (Cimpoi et al.,
 955 2014), Places365 (López-Cifuentes et al., 2020)
- 957 • **CIFAR-100**
 Near-OOD: CIFAR-10 (Krizhevsky et al., 2009), Tiny ImageNet (Le & Yang, 2015)
 Far-OOD: MNIST (LeCun et al., 1998), SVHN (Netzer et al., 2011), Textures (Cimpoi et al.,
 959 2014), Places365 (López-Cifuentes et al., 2020)
- 961 • **ImageNet**
 Near-OOD: SSB-hard (Vaze et al., 2021), NINCO (Bitterwolf et al., 2023)
 Far-OOD: iNaturalist (Van Horn et al., 2018), Textures (Cimpoi et al., 2014), OpenImage-O
 (Wang et al., 2022)

965 B.3 ABLATION STUDY

966
 967 In this section, we highlight the relevant parameters for the ablation study experiments on SISOM.
 968 Namely, we examine the Optimal Sigmoid Steepness (OS) and the Reduced Subset Selection (RS)
 969 shown in Tab. 4. In the experiments conducted with RS, a representative subset size of 10% relative
 970 to the original training set was used across all experiments. Additionally, the specific distance radius
 971 used for the class-wise ProbCover (Yehuda et al., 2022) implementation on CIFAR-10, CIFAR-100,
 and ImageNet is provided in Table 4. For SISOM + RS without OS, the suggested sigmoid values

Table 4: Parameters for the Ablation Study, Probcover Radius for RS and Search Space of Optimal Sigmoid Steepness.

Dataset	ProbCover Radius	Layer	Sigmoid Search Values
CIFAR-10	0.75	AdaptiveAvgPool2d-1	100 , 1000
		Sequential-3	1 , 10, 1000
		Sequential-2	0.001 , 0.1, 1
		Sequential-1	0.001 , 0.1, 1
CIFAR-100	5.0	AdaptiveAvgPool2d-1	1 , 50, 100
		Sequential-3	0.1 , 10, 100
		Sequential-2	0.1 , 1
		Sequential-1	0.005, 0.1
ImageNet	10.0	AdaptiveAvgPool2d-1	10, 100, 3000
		Sequential-3	1 , 10, 300
		Sequential-2	0.1 , 1
		Sequential-1	0.1 , 1

(Liu et al., 2024) emphasized in Appendix B.2 were used. For the OS modification, the search space for the optimal sigmoid parameters is presented in Table 4. The parameters fulfilling the minimization of Eq. (10) are highlighted in bold.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025