

# Active Learning for Object Detection with Non-Redundant Informative Sampling

Aral Hekimoglu  
Technical University Munich  
Munich, Germany  
aral.hekimoglu@tum.de

Adrian Brucker  
Technical University Munich  
Munich, Germany  
adrian.brucker@tum.de

Alper Kagan Kayali  
Technical University Munich  
Munich, Germany  
alper.kagan.kayali@tum.de

Michael Schmidt  
BMW Group  
Munich, Germany  
michael.se.schmidt@bmw.de

Alvaro Marcos-Ramiro  
BMW Group  
Munich, Germany  
alvaro.marcos-ramiro@bmw.de

## Abstract

Curating an informative and representative dataset is essential for enhancing the performance of 2D object detectors. We present a novel active learning sampling strategy that addresses both the informativeness and diversity of the selections. Our strategy integrates uncertainty and diversity-based selection principles into a joint selection objective by measuring the collective information score of the selected samples. Specifically, our proposed NORIS algorithm quantifies the impact of training with a sample on the informativeness of other similar samples. By exclusively selecting samples that are simultaneously informative and distant from other highly informative samples, we effectively avoid redundancy while maintaining a high level of informativeness. Moreover, instead of utilizing whole image features to calculate distances between samples, we leverage features extracted from detected object regions within images to define object features. This allows us to construct a dataset encompassing diverse object types, shapes, and angles. Extensive experiments on object detection and image classification tasks demonstrate the effectiveness of our strategy over the state-of-the-art baselines. Specifically, our selection strategy achieves a 20% and 30% reduction in labeling costs compared to random selection for PASCAL-VOC and KITTI, respectively.

## 1. Introduction

Accurately detecting 2D objects [9, 24, 36] is fundamental in scene understanding across various applications, such as autonomous driving. The current state-of-the-art (SOTA) object detection systems rely on deep learning techniques

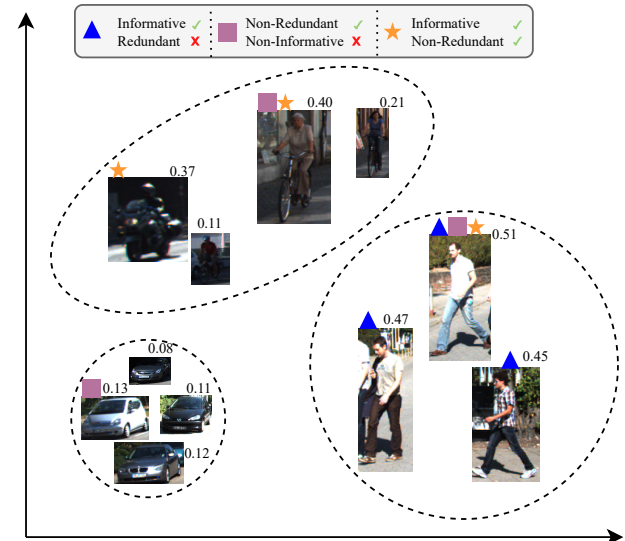


Figure 1. Illustration of different AL selection strategies in feature space. The size of each object symbolizes its uncertainty score, also denoted next to the object. Dashed lines mark clusters formed by k-means. The uncertainty-based method [37] (triangle) selects objects with the highest uncertainty. However, due to the similarity among the selected objects, there is redundancy in the selections. The diversity-based method [29] (square) selects one object from each cluster. However, since the selected *Car* belongs to a well-represented class and has low uncertainty, its selection is less informative. Our proposed approach, NORIS (star) selects object with high uncertainty while maintaining diversity.

and heavily depend on large-scale annotated datasets. However, annotating extensive amounts of 2D detection data is time-consuming and labor-intensive, making labeling the

entirety of the collected data challenging. Moreover, selecting similar samples might lead to overfitting, where the model becomes overly specialized in some areas of the input space, limiting its generalization ability. Therefore, curating an informative dataset that is also representative of the target distribution is crucial to achieving optimal object detection performance.

Active learning (AL) [11, 14, 29] emerged as a promising solution for selecting the most suitable data for labeling. Uncertainty-based AL methods [7, 8, 18, 23, 27, 34] leverage the uncertainty of the output of the network to assign an information score to each sample and subsequently select samples with the highest score for labeling. However, as depicted in Fig. 1, these methods tend to select highly similar samples, such as the three pedestrians with similar poses, leading to redundancy and a lack of diversity within the labeled dataset.

To prevent redundancy in selections, diversity-based methods [1, 29, 31, 32] select representative samples that cover the input space. One approach is to cluster the data and choose one sample from each cluster. However, selecting the most diverse samples does not necessarily guarantee informativeness. As shown in Fig. 1, the diversity-based method selects a *Car* because it belongs to a distinct cluster from the other selected samples. However, this selection is less informative as the car has a lower uncertainty compared to other samples.

To address the limitations of the aforementioned methods, hybrid approaches [5, 20, 30, 33] integrate uncertainty and diversity by multiplying the individual uncertainty and diversity scores. However, these methods treat uncertainty and diversity as separate properties and do not fully capture the interaction between sample uncertainties. For instance, in Fig. 1, learning from one of the pedestrians would decrease the uncertainty of nearby pedestrians. To better estimate the collective information score of a selected set, modeling the interaction of uncertainties based on the distance of samples is necessary.

We propose a novel sampling strategy called Non-Redundant and Informative Sampling (NORIS) in Sec. 3.2. NORIS tackles the challenges of redundancy and informativeness by formulating the active learning objective as a collective information score. Our work is the first to formally integrate the interaction (as defined by the diversity) of different uncertainties within the same AL selection objective and offer an effective algorithmic solution.

Furthermore, previous AL methods for object detection primarily rely on the distance between image features to define diversity. As a result, these methods select a diverse range of scenes, including variations in lighting and weather conditions. However, their holistic image approach makes it difficult to capture the distinct diversity of objects within these scenes. To overcome this limitation, in

Sec. 3.3, we propose to calculate the distance using object features, which we define as the intermediate network features within the region of interest (ROI) for each detected box. This allows us to capture the diversity of objects within a scene. With our approach, the detector is trained on a more comprehensive and diverse set of object instances, enabling it to generalize better and achieve improved detection performance across different scenarios.

Our main contributions are summarized as follows:

- We propose NORIS, a novel strategy that addresses the challenge of selecting informative samples while minimizing redundancy by formulating the AL selection objective as a collective information score.
- We introduce a unique approach to define diversity by utilizing object features extracted from the detector, enhancing the generalizability of the detector to a diverse set of objects.
- We conduct an extensive experimental evaluation to assess the performance of NORIS against the existing SOTA AL baselines, demonstrating superior performance and up to 30% reduction in labeling costs.

## 2. Related work

### 2.1. Uncertainty-based methods

Most AL methods for object detection rely on an information score, such as the uncertainty of the model in its predictions, to prioritize the selection of the most informative samples for labeling. For object detection, uncertainty-based methods typically target either the classification uncertainty or the localization uncertainty of the predicted bounding box. Classification uncertainty-focused methods utilize the predicted class probability distribution as an estimate of the model’s confidence for each sample [2, 4, 16, 27, 39]. Localization uncertainty methods [7, 19, 28] use the uncertainty associated with the bounding box parameters such as position or dimensions to formulate an information score. Kao *et al.* [19] introduced two uncertainty measures, namely localization tightness and localization consistency, for sample selection. Choi *et al.* [7] proposed an aleatoric uncertainty head designed to estimate localization uncertainty.

Another approach to quantify informativeness involves measuring the disagreement among different predictions for the same input. Epistemic uncertainty-based methods leverage diverse predictions obtained from multiple models, where each model is trained with different initial random weights [6], or from a single model with a dropout layer that generates distinct predictions at each forward pass [25]. In inconsistency-based methods [11, 38], diverse predictions are obtained by applying various augmentations to the same input sample.

A task-agnostic AL approach, LL4AL, was proposed by Yoo *et al.* [37]. They include an additional head in the network designed to predict the target loss for unlabeled inputs, and the information score is defined based on this predicted loss.

## 2.2. Diversity-based active learning

A common issue with uncertainty-based methods is the tendency to select similar samples, leading to redundant selections. In contrast, diversity-based methods [10,29,32,35] select representative samples that cover the entire dataset space. These methods apply a similarity measure between a candidate and previously selected samples, choosing the sample with the least similarity to the previous selections. Diversity-based methods can be categorized into three groups based on the type of similarity measurement: feature-based [29], adversarial [10,32], and context-based [1].

Feature-based methods represent each sample with a feature vector and employ a distance metric to estimate the similarity. Sener *et al.* [29] proposed constructing a core-set by efficiently solving the k-center problem, where distances are computed as the Euclidean norm between features extracted from intermediate layers of a network for each image. Adversarial diversity-based methods [32] train a discriminator to predict whether a sample belongs to the previously selected set. The output from the discriminator is then used as the similarity measure to the labeled subset. Agarwal *et al.* [1] introduced a method for contextual diversity, in which the KL-divergence between the predicted class probabilities of a sample and the previously selected subset is used to define contextual similarity.

Existing diversity-based methods primarily focus on images and have not been adapted to object detection. Haussmann *et al.* [16] implemented a feature-based diversity method for object detection but relied on image features. In our work, we propose leveraging object features extracted from intermediate layers of an object detector, aiming to construct a training set with a more diverse set of objects.

## 2.3. Hybrid uncertainty-diversity methods

Considering both uncertainty and diversity is critical in forming an effective selection strategy that identifies informative samples while simultaneously avoiding redundancy [3,5,20,30,31,33,40]. Shen *et al.* [31] proposed K-Covers, a method that first clusters the unlabeled set using a modified k-means algorithm and then selects the most uncertain sample from each cluster. Zhdanov *et al.* [40] proposed DBAL, which leverages a weighted k-means clustering algorithm by multiplying the distance by the uncertainty to prioritize both diversity and informativeness of a sample. Haussmann *et al.* [16] adopt a similar approach by multiplying the distance between feature vectors with the uncertainty and se-

lecting samples based on the resultant scores. Ash *et al.* introduced BADGE [3], which uses gradient embeddings as a proxy for uncertainty and iteratively selects a batch of samples that are both diverse and uncertain based on their distances in the gradient embedding space.

In contrast to prior hybrid methods that derive a singular selection score by multiplying the uncertainty and diversity of a sample, our approach focuses on the information score of a set of samples considered collectively. We propose a unique approach to integrate uncertainty and diversity by assessing the impact that labeling and training a sample would have on the information score of other samples, with the impact scaled based on their similarities.

## 3. Methodology

### 3.1. Problem definition

Let  $(x, y)$  denote a sample pair drawn from a labeled training dataset  $(X_{\text{train}}, Y_{\text{train}})$  where  $X_{\text{train}}$  represents the set of labeled data points and  $Y_{\text{train}}$  their corresponding labels. Let  $u$  denote an unlabeled sample drawn from a larger pool of unlabeled samples  $X_U$ . In each AL cycle, a model  $\phi$  is trained with the labeled data  $(X_{\text{train}}, Y_{\text{train}})$ . Then, an acquisition function selects a subset of unlabeled samples to be manually labeled by an external oracle. These newly labeled samples are then added to the labeled set for training in the next iteration. We denote the selected subset by  $S \subseteq X_U$ , and  $B = |S|$  represents the labeling budget in each iteration. Consequently, the AL objective can be defined as selecting a subset  $S$  of unlabeled data points such that the performance of the model trained on  $(X_{\text{train}} \cup S, Y_{\text{train}} \cup Y_S)$  is maximized, where  $Y_S$  represents the labels for the samples in  $S$ .

### 3.2. Non-redundant informative sampling

NORIS encapsulates proposed query strategies that combine uncertainty with diversity by incorporating redundancy. We define a function  $\sigma : X \rightarrow \mathbb{R}_{\geq 0}$  that assigns each unlabeled data  $u \in X_U$  an information score  $\sigma(u)$  (e.g., the model’s uncertainty for that image). In the setting of uncertainty-based AL, in each query iteration, the  $B$  highest-ranked samples are selected for annotation. This can be formulated as maximizing the sum of their information scores:

$$\arg \max_{S \subseteq X_U, |S|=B} \sum_{u \in S} \sigma(u) \quad (1)$$

One drawback of this approach is that the information gain from a sample is reduced if we already select and train with a similar, more informative sample.

To model this interaction, we consider the loss of a sample as our selection criteria. Training on one sample reduces the loss of similar samples as the network learns to generalize in that region in the optimization space. Formally, we

follow the proof in Core-Set [29] and define the loss function to be  $\kappa$ -Lipschitz continuous. With this definition, we express the upper bound on the loss of one sample  $u$  after training with another sample  $v$  (details provided in the supplementary material):

$$l(u; \theta_{i+1}) \leq l(u; \theta_i) + 2 \cdot \kappa \cdot d(u, v) - l(v; \theta_i) \quad (2)$$

where  $d(u, v)$  represents the distance between  $u$  and  $v$ , and  $\theta_i$  and  $\theta_{i+1}$  denote the parameters of the network before and after a training iteration that includes sample  $v$ . Following this argumentation, the final loss of a sample  $u$  is bounded by an expression inversely related to the initial loss of sample  $v$  and directly related to the distance between samples  $u$  and  $v$ . This implies that closer samples have a stronger influence on reducing each other’s loss, while higher initial loss of the trained sample  $v$  leads to a tighter bound in the loss of the neighboring sample  $u$ .

Therefore, to get an information score of a set collectively, we need to model the interaction of training on one sample would have on the information score of other samples. We propose a model to update the information score of the sample  $u$ , assuming that  $v$  is selected.

$$\sigma'(u) := \sigma(u) - \text{sim}(u, v) * \sigma(v) \quad (3)$$

Here, we define a function that measures the similarity between samples  $\text{sim}(u, v)$  to be inversely correlated to  $d(u, v)$ . We require that  $\text{sim}(u, u) = 1$  for all  $u \in X_U$  and that  $\text{sim}$  is symmetric, i.e.  $\text{sim}(u, v) = \text{sim}(v, u)$  for all  $u, v \in X_U$ . The more similar  $u$  and  $v$  are, the higher  $\text{sim}(u, v)$  should be. We propose two similarity measures that utilize the distance metric  $d$ :

$$\text{Gaussian: } \text{sim}(u, v) = e^{-\frac{1}{\lambda} d(u, v)^2} \quad (4)$$

$$\text{linear: } \text{sim}(u, v) = \max \left\{ 0, 1 - \frac{d(u, v)}{\lambda} \right\} \quad (5)$$

The hyperparameter  $\lambda > 0$  indicates the influence of a sample to its surroundings. Higher values represent stronger influence, while lower values lead to weaker influence. An exhaustive hyperparameter search can assist in finding an appropriate value. For the linear similarity, we suggest experimenting with values in the range of  $(0, d_{\max}]$ , where  $d_{\max}$  is defined as:

$$d_{\max} := \max_{x, x' \in X_U} d(x, x'). \quad (6)$$

We argue that values for  $\lambda$  smaller than or equal to the maximum distance between the unlabeled samples are reasonable, as the similarity score of the most distant samples in the embedding space should be zero. Note that  $d_{\max}$  must be recalculated in each AL cycle since distances are calculated between the changing low-dimensional embeddings.

This also allows for automatic adaptation to the scale of the feature space as opposed to fixing a value throughout all AL iterations. When conducting a hyperparameter search for the linear similarity, we fix a value  $\alpha \in (0, 1]$  and scale it to  $\lambda = \alpha \cdot d_{\max}$  in each AL cycle. For the Gaussian similarity, we fix a value  $\alpha \in (0, 1]$  and scale it to  $\lambda = \frac{(\alpha \cdot d_{\max})^2}{\pi}$  in each AL cycle for hyperparameter tuning. For additional details, please refer to the supplementary material.

We propose two variants on how to incorporate Eq. (3) into an AL selection objective called NORIS-Sum and NORIS-Max. Both versions are not restricted to a specific method of measuring the information value of individual images. The choice of similarity measure between images is also flexible, making NORIS highly versatile.

**NORIS-Sum.** In NORIS-Sum, we define the AL selection objective as the aggregate of the information scores. We determine the individual score of a sample  $u$  as the information score of  $u$ , deducted by the weighted sum of information scores of all the selected samples  $v \in S \setminus \{u\}$ , scaled by their similarity  $\text{sim}(u, v)$ . We sum the individual scores to generate a batch-wise acquisition score that better captures the true value of a set of samples.

$$\arg \max_{S \subseteq X_U, |S|=B} \sum_{u \in S} \sigma(u) - \left( \sum_{v \in S \setminus u} \text{sim}(u, v) * \sigma(v) \right) \quad (7)$$

To efficiently solve this objective, we propose the NORIS-Sum algorithm (Algorithm 1). In each iteration of the while-loop, the sample  $u_{\text{next}}$  that leads to the highest immediate gain is selected and added to  $S$ . Subsequently, for each  $u \in X_U \setminus S$ , the information score is updated by subtracting the uncertainty of  $u_{\text{next}}$  scaled by the similarity  $\text{sim}(u, u_{\text{next}})$ . This algorithm has a runtime complexity of  $\mathcal{O}(B \cdot |X_U|)$ .

---

#### Algorithm 1 NORIS-Sum

---

**Require:** unlabeled dataset  $X_U$ , batch size  $1 \leq B \leq |X_U|$

```

 $S \leftarrow \emptyset$ 
while  $|S| \neq B$  do
     $u_{\text{next}} \leftarrow \arg \max_{u \in X_U \setminus S} \sigma(u)$ 
     $S \leftarrow S \cup \{u_{\text{next}}\}$ 
    for  $u \in X_U \setminus S$  do
         $\sigma(u) \leftarrow \sigma(u) - \text{sim}(u, u_{\text{next}}) * \sigma(u_{\text{next}})$ 
    end for
end while

```

---

**NORIS-Max.** In NORIS-Max, instead of reducing the information score of an individual image for each sample in its vicinity, we consider only the impact of the closest one. We argue that distant data has less influence, as the loss bound in Eq. (2) becomes less tight as the distance increases. Accumulating many small similarity scores can



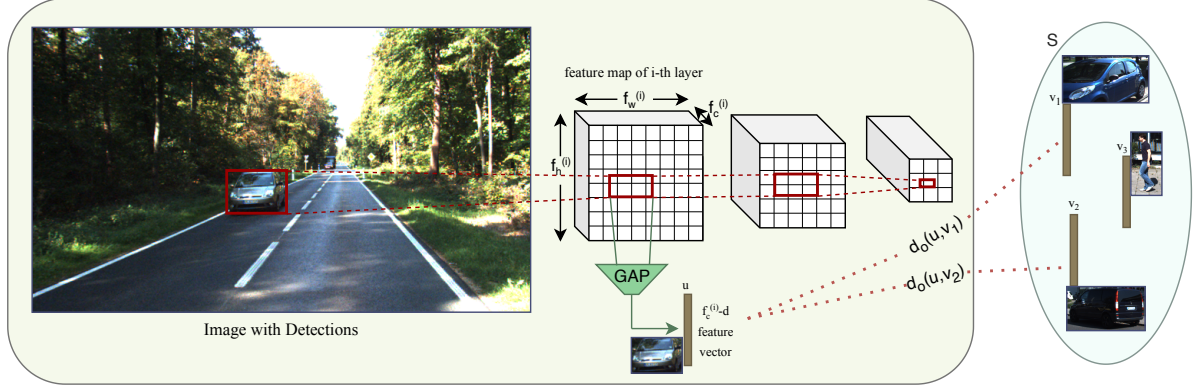


Figure 2. Illustration of the object-diversity calculation. For each detected object, we crop the corresponding ROI from intermediate feature maps and apply global average pooling to get a feature vector. We then measure the distance between the feature vector of the detected object and the feature vectors of objects in selected set  $S$  to define the distance metric.

have a strong impact which might not be desirable. Therefore, in this approach, the score of a sample is reduced only by the similarity of its closest sample in the batch of selected points. The batch-wise acquisition function becomes:

$$\arg \max_{S \subseteq X_U, |S|=B} \sum_{u \in S} \sigma(u) - \text{sim}(u, v_c) \cdot \sigma(v_c) \quad (8)$$

where  $v_c$  represents the most similar sample to sample  $u$ . We follow the same algorithm described in Algorithm 1, and modify the update step to only consider its closest sample.

### 3.3. Object diversity

One of the critical design choices in our algorithm is the distance metric  $d(u, v)$ . While the Euclidean distance between intermediate features of a deep neural network is a common choice for image classification [29], we propose a novel approach that uses features describing the object instead of features for the whole image. This approach is more intuitive and aligns better with the goal of covering a diverse set of objects with different variations in types, angle, and occlusion patterns to make the object detector more generalizable.

Our approach, illustrated in Fig. 2, leverages object features to compute the similarity between objects. We use our object detector to detect all objects in the unlabeled pool and extract intermediate feature maps  $f^{(i)}$  with dimensions of  $(f_w^{(i)}, f_h^{(i)}, f_c^{(i)})$ . We compute the ROI for each detected object in the feature map  $f^{(i)}$ , defined by its bounding box coordinates. We scale the location  $(o_x * f_w^{(i)} / i_w, o_y * f_h^{(i)} / i_h)$  and size  $(o_w * f_w^{(i)} / i_w, o_h * f_h^{(i)} / i_h)$ , where  $(o_x, o_y, o_w, o_h)$  represent the pixel-wise location, width, and height of the detected object, and  $(i_w, i_h)$  represent the width and height of the input image. Next, we crop the corresponding region from the feature map, apply

global average pooling (GAP), and obtain a feature vector with  $f_c^{(i)}$  elements. These features are then used to compute the Euclidean distance between objects.

To compute the distance between images based on the objects they contain, we aggregate the distances between the objects by taking the maximum distance between any two objects in the images. Thus, for two images  $u$  and  $v$ , the object-based distance  $d_o(u, v)$  is defined as follows:

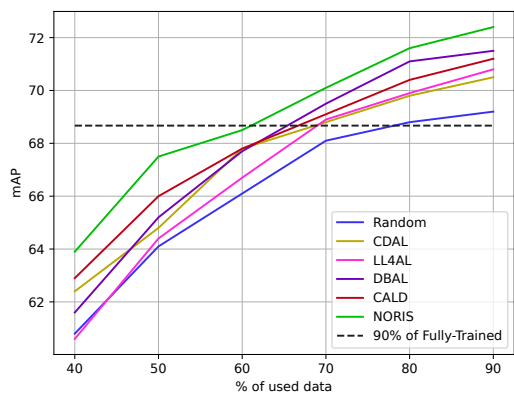
$$d_o(u, v) = \max_{o_u \in u, o_v \in v} |f_{o_u} - f_{o_v}|^2 \quad (9)$$

Here,  $o_u$  and  $o_v$  represent the sets of detected objects in the images  $u$  and  $v$ , respectively, and  $f_{o_u}$  and  $f_{o_v}$  represent the feature vectors for each object. The maximum distance is selected since the overall distance between images should reflect the distance between the most dissimilar objects in the images.

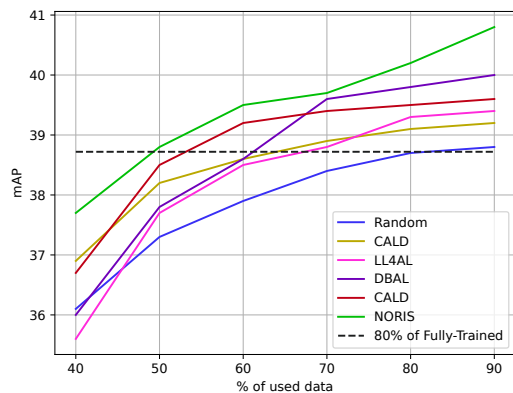
In addition to object features, we incorporate scene diversity into our method to capture a diverse set of scenes with varying weather and lighting conditions and to be able to measure the distance between two images in the absence of detected objects. To obtain image features, we adopt a similar approach as with object features by applying global average pooling with ROI defined as the entire image. For any two images  $u$  and  $v$ , we define the final distance  $d(u, v)$  as:

$$d(u, v) = d_o(u, v) * |f_u - f_v|^2 \quad (10)$$

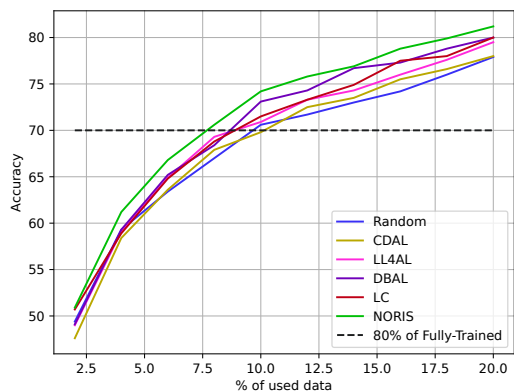
Here,  $f_u$  and  $f_v$  are the global image features for images  $u$  and  $v$ , respectively. This definition of distance captures both object and scene diversity, allowing our algorithm to select samples that cover a diverse range of object variations and scenes.



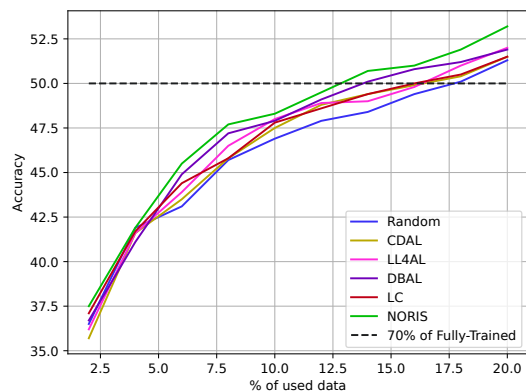
(a) PASCAL-VOC 2007



(b) KITTI val



(c) CIFAR 10



(d) CIFAR 100

Figure 3. Comparison with SOTA AL methods. Lines indicate the averaged results over three trials. Note that all methods start from the same network trained with the initially labeled data.

## 4. Experiments

### 4.1. Comparison with baselines

**Datasets and evaluation metric.** We evaluate the performance of our proposed algorithm on four datasets, two for object detection: PASCAL VOC [12], and KITTI [15], and two for image classification: CIFAR-10 [21], and CIFAR-100 [21]. We use the mean Average Precision (mAP) [13] as the evaluation metric for object detection and accuracy for image classification. For PASCAL VOC, we combine the 2007 and 2012 datasets, resulting in 16,551 samples, which serve as the unlabeled pool. We evaluate the performance on the test set of PASCAL VOC 2007. For KITTI, we follow the data split used in [22]. The training set consists of 7,481 images, which are divided into an unlabeled pool with 3,712 samples and a validation set with 3,769 samples. We report the results on the validation set. For the CIFAR-10 and CIFAR-100 datasets, we

use the training set containing 50,000 images as the initial unlabeled pool and report results on the test set.

**Model architecture.** We use CenterNet [41] as our base detector. We train the model for the same number of epochs and follow the same hyperparameters and optimization scheme described in the original paper. The experiments are conducted on an NVIDIA Tesla V100 GPU using the PyTorch deep learning framework [26]. For the image classification experiments, we use ResNet-18 [17] as our base classifier and extract the image feature embeddings from the penultimate layer to calculate distances. In each query iteration, we train the model continuously for 50 epochs with an early stopping patience of 10, using the Adam optimizer with a learning rate of  $10^{-3}$ .

**Active learning details.** We start our experiments with an unlabeled pool of samples and create an initial training set by randomly selecting samples. We randomly split the unlabeled pool into a labeled pool, which serves as the ini-

tial training set. At each AL iteration, we select samples from the remaining unlabeled pool to add to the training set using the selection strategy. To simulate the labeling process, we use the already available annotations. For the object detection datasets, we initially select 30% of the available samples and add 10% more during each AL cycle. For image classification tasks, we initially select 1000 images (equivalent to 2% of the total dataset) and add 1000 images in each iteration. We report the mean of the respective metric for each experiment over three independent runs with different random selections of the initial labeled pool. We also report the variances across these three runs in the supplementary material.

**Baselines.** To evaluate the effectiveness of our proposed selection strategy, we compare it against several baselines from the literature. We choose two SOTA uncertainty-based methods: the inconsistency-based selection strategy CALD [38], and the loss-based method LL4AL [37]. Additionally, we compare against the diversity-based method CDAL [1]. For hybrid methods, we include DBAL [40], which ranks samples based on the multiplication of their diversity and uncertainty scores. To mimic passive learning, we use **Random** selection, where each sample is assigned a score following a uniform distribution. To define the uncertainty function in NORIS, we use the least-confidence approach, where the uncertainty is determined by the predicted class probability. To provide a baseline for comparison, we report the results obtained from a "fully-trained" network trained on the entire training set.

**Results on PASCAL VOC.** Our experimental results on the PASCAL VOC dataset, presented in Fig. 3a, demonstrate the performance of our proposed selection strategy over the baseline methods. Our method outperforms all baselines by at least 1.0 mAP in the initial AL cycle. In the second cycle, our method outperforms the Random baseline by 3.4 mAP and the second-best method, CALD, by 1.5 mAP. As the number of actively selected samples grows, the two methods that combine uncertainty and diversity, NORIS and DBAL, consistently outperform the other baselines, highlighting the significance of utilizing uncertainty and diversity jointly in active learning. In the final cycle, where we train our detector with 90% of the data, of which 60% are actively sampled using our selection strategy, our proposed algorithm outperforms all baseline methods by at least 0.9 mAP. Notably, our approach achieves 90% of the performance of a fully-trained model with only 60% of the samples. This corresponds to 20% less labeled data compared to the Random baseline and 5% less labeled data compared to CALD. These results demonstrate that our proposed algorithm effectively selects informative samples while minimizing redundancy in the labeled dataset, resulting in improved data-efficiency and accuracy.

**Results on KITTI.** The performance of our proposed al-

Variant	Similarity	mAP
NORIS-Max	linear	40.1
	gaussian	<u>40.8</u>
NORIS-Sum	linear	40.4
	gaussian	<b>41.0</b>

Table 1. Ablation study on variants of NORIS and similarity measures. Bold indicates the best-performing setting, and underline indicates the best-performing similarity measure within the NORIS-Max variant.

gorithm and the baseline methods on the KITTI dataset are presented in Fig. 3b. In the initial AL cycle, our method outperforms the Random baseline by 1.6 mAP and achieves a 1.0 mAP lead over the second-best method, CALD. Our approach maintains its lead over the baselines in subsequent cycles, and in the final cycle, where 60% of the samples are actively selected, we outperform all the baselines by 0.8 mAP. Our proposed algorithm achieves 80% of the fully-trained performance using only 50% of the data, corresponding to a 30% improvement in data savings compared to the Random baseline’s 80%.

**Results on CIFAR.** In CIFAR-10 (Fig. 3c), our method reaches the 80% of the fully-trained performance using only 7.5% of the data, with a data savings rate of 1.0% compared to DBAL and 2.5% compared to Random. In the final cycle, where we train with 20% of the available data, our method reaches 1.0 higher accuracy. In CIFAR-100 (Fig. 3d), our method achieves 70% of the fully-trained performance using only 12.5% of the data, with a data savings rate of 1.0% compared to DBAL and 5.0% compared to Random. In the final cycle, where we train with 20% of the available data, our method reaches 1.3 higher accuracy.

## 4.2. Ablation on NORIS

**Variants of NORIS.** As we proposed two variants of NORIS - namely NORIS-Sum and NORIS-Max - along with two types of similarity scores (linear and Gaussian) in Sec. 3.2, we perform an ablation study to determine their relative performance. Our results, presented in Tab. 1, indicate that NORIS-Sum tends to outperform NORIS-Max, and the Gaussian similarity score leads to superior performance over the linear similarity. We follow this configuration throughout our experiments.

**t-SNE of selections.** To offer a visual perspective on how different selection methods operate, we provide a t-SNE visualization of selected points. The uncertainty-based method CALD, as shown in Fig. 4a, tends to select highly uncertain samples. However, it disregards samples situated in the outer regions, leading to the selection of numerous close-proximity samples, which introduced redundancy. CDAL, a diversity-based method illustrated in Fig. 4b, selects more samples from the outer regions but tends to over-

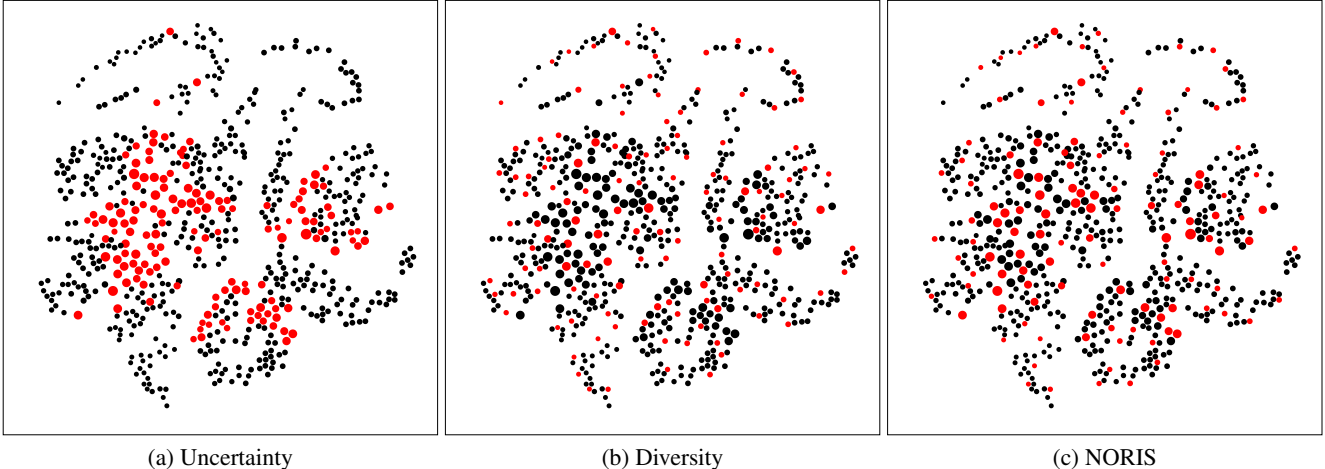


Figure 4. t-SNE visualizations of feature embeddings for the CIFAR-10 dataset. Each dot corresponds to an image embedding, with its size being scaled in accordance with the uncertainty of the image. The red dots indicate the samples selected for labeling corresponding to the respective selection strategy.

look some highly uncertain objects. NORIS, visualized in Fig. 4c, strikes a balance between selecting samples that are both informative and from outer regions.

### 4.3. Ablation on Object Features

**Comparison of using image or object features.** To demonstrate the effectiveness of our proposal to use object features in contrast to image features for defining diversity, we adapt three SOTA diversity-based methods to work with object features: Core-Set [29] VAAL [32] and MAL [10]. Specifically, we use features from the object detector and modify the discriminator heads of VAAL and MAL to output a similarity value between 0 and 1 for each detected object. We then compute a diversity score for each image by assigning the score of the discriminator output of its most dissimilar object. We refer the reader to the supplementary material for more detailed information on the modifications and a comparison of the layers from which the features are extracted within the CenterNet architecture.

We present our comparison results on the KITTI dataset in Tab. 2. Object features consistently outperform image features in all cycles for the MAL method. For Core-Set and VAAL, using image features initially leads to better performance, but using object features results in superior performance in later cycles. This observation can be attributed to the lower performance of the object detector during the initial cycles, making it more challenging to rely on the detected objects for diversity. Overall, the findings indicate that object features improve the performance of diversity-based AL algorithms compared to image features.

**Visual comparison.** In Fig. 5, we compare qualitatively the selection strategy of Core-Set with object features. Our analysis reveals that the query objects exhibit similar char-

		40%	50%	70%	90%
Core-set [29]	Img	<b>36.9</b>	<b>38.2</b>	38.5	38.7
	Obj	36.8	38.0	<b>38.6</b>	<b>38.9</b>
VAAL [32]	Img	<b>36.6</b>	<b>38.2</b>	38.5	39.1
	Obj	36.3	38.1	<b>39.0</b>	<b>39.4</b>
MAL [10]	Img	36.7	37.9	38.8	39.2
	Obj	<b>37.1</b>	<b>38.3</b>	<b>39.2</b>	<b>39.6</b>

Table 2. Comparison of using image or object features for three different diversity-based methods on KITTI. Bold values indicate the best performance for each method.

acteristics to their most similar object. Objects belong to the same class and have similar visual appearances, like the color of the vehicles in the second row. This suggests that using object features can identify these patterns and select dissimilar objects, resulting in a more diverse and representative training set.

## 5. Conclusion

We introduced NORIS, a novel AL strategy that combines uncertainty with diversity to select the most informative and non-redundant samples for object detection. Our approach jointly optimizes both criteria by reducing the information score of a sample if it is close to other selected samples with high uncertainty. To better define diversity for object detection, we proposed using the distance between object features to measure the similarity of samples. Our experiments demonstrated that incorporating non-redundant sampling with object feature-based diversity led to better performance than SOTA AL methods. Moreover, we observed improved performance for all diversity-based meth-



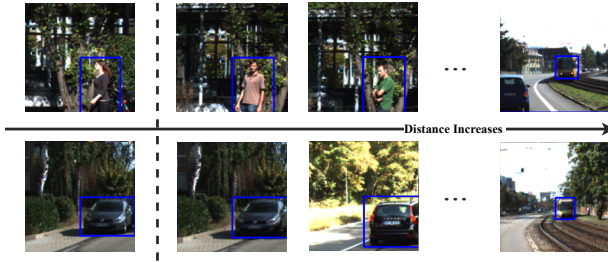


Figure 5. Visualization of the distances computed using object features. The leftmost column corresponds to the selected query objects. For each query object, the available objects from the unlabeled dataset are ranked in ascending order of their distance to the respective query sample from left to right.

ods when utilizing object features.

## References

- [1] Sharat Agarwal, Himanshu Arora, Saket Anand, and Chetan Arora. Contextual diversity for active learning. In *ECCV*, 2020. 2, 3, 7
- [2] Hamed H Aghdam, Abel Gonzalez-Garcia, Joost van de Weijer, and Antonio M López. Active learning for deep detection neural networks. In *ICCV*, 2019. 2
- [3] Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep batch active learning by diverse, uncertain gradient lower bounds. In *ICLR*, 2019. 3
- [4] Clemens-Alexander Brust, Christoph Käding, and Joachim Denzler. Active learning for deep object detection. In *VISAPP*, 2019. 2
- [5] Xiaofeng Cao and Ivor W. Tsang. Bayesian active learning by disagreements: A geometric perspective. *CoRR*, 2021. 2, 3
- [6] Kashyap Chitta, José M Álvarez, Elmar Haussmann, and Clément Farabet. Training data subset search with ensemble active learning. *T-ITS*, 2021. 2
- [7] Jiwoong Choi, Ismail Elezi, Hyuk-Jae Lee, Clément Farabet, and José Manuel Álvarez. Active learning for deep object detection via probabilistic modeling. In *ICCV*, 2021. 2
- [8] Sai Vikas Desai, Akshay L Chandra, Wei Guo, Seishi Nishimiyama, and Vineeth N. Balasubramanian. An adaptive supervision framework for active learning in object detection. In *BMVC*, 2019. 2
- [9] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. Centernet: Keypoint triplets for object detection. In *ICCV*, 2019. 1
- [10] Sayna Ebrahimi, William Gan, Dian Chen, Giscard Biamby, Kamyar Salahi, Michael Laielli, Shizhan Zhu, and Trevor Darrell. Minimax active learning. *arXiv preprint arXiv:2012.10467*, 2020. 3, 8, 12
- [11] Ismail Elezi, Zhiding Yu, Anima Anandkumar, Laura Leal-Taixe, and Jose M Alvarez. Not all labels are equal: Rationalizing the labeling costs for training object detection. In *CVPR*, 2022. 2
- [12] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 2010. 6
- [13] Di Feng, Ali Harakeh, Steven L Waslander, and Klaus Dietmayer. A review and comparative study on probabilistic object detection in autonomous driving. *T-ITS*, 2021. 6
- [14] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In *ICML*, 2017. 2
- [15] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, 2012. 6
- [16] Elmar Haussmann, Michele Fenzi, Kashyap Chitta, Jan Ivanek, Hanson Xu, Donna Roy, Akshita Mittel, Nicolas Koumchatzky, Clement Farabet, and Jose M Alvarez. Scalable active learning for object detection. In *IV*, 2020. 2, 3
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 6
- [18] Aral Hekimoglu, Michael Schmidt, Alvaro Marcos-Ramiro, and Gerhard Rigoll. Efficient active learning strategies for monocular 3d object detection. In *IV*, 2022. 2
- [19] Chieh-Chi Kao, Teng-Yok Lee, Pradeep Sen, and Ming-Yu Liu. Localization-aware active learning for object detection. In *ACCV*, 2018. 2
- [20] Andreas Kirsch, Joost van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning. *CoRR*, 2019. 2, 3
- [21] Alex Krizhevsky. Learning multiple layers of features from tiny images. *CiteSeer*, 2009. 6
- [22] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *CVPR*, 2019. 6
- [23] Ying Li, Binbin Fan, Weiping Zhang, Weiping Ding, and Jianwei Yin. Deep active learning for object detection. *Information Sciences*, 2021. 2
- [24] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. In *CVPR*, 2022. 1
- [25] Dimity Miller, Lachlan Nicholson, Feras Dayoub, and Niko Sünderhauf. Dropout sampling for robust object detection in open-set conditions. In *ICRA*, 2018. 2
- [26] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimeshain, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, 2019. 6
- [27] Soumya Roy, Asim Unmesh, and Vinay P Nambodiri. Deep active learning for object detection. In *BMVC*, 2018. 2
- [28] Sebastian Schmidt, Qing Rao, Julian Tatsch, and Alois Knoll. Advanced active learning strategies for object detection. In *IV*, 2020. 2
- [29] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. In *ICLR*, 2018. 1, 2, 3, 4, 5, 8, 11
- [30] Dhruv Sharma, Zahil Shanis, Chandan K Reddy, Samuel Gerber, and Andinet Enquobahrie. Active learning technique

- for multimodal brain tumor segmentation using limited labeled images. In *MICCAIW*, 2019. 2, 3
- [31] Yeji Shen, Yuhang Song, Hanhan Li, Shahab Kamali, Bin Wang, and C-C Jay Kuo. K-covers for active learning in image classification. In *ICMEW*, 2019. 2, 3
- [32] Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell. Variational adversarial active learning. In *ICCV*, 2019. 2, 3, 8, 12
- [33] Yao Tan, Liu Yang, Qinghua Hu, and Zhibin Du. Batch mode active learning for semantic segmentation based on multi-clue sample selection. In *CIKM*, 2019. 2, 3
- [34] Fuhui Tang, Dafeng Wei, Chenhan Jiang, Hang Xu, Andi Zhang, Wei Zhang, Hongtao Lu, and Chunjing Xu. Towards dynamic and scalable active learning with neural architecture adaption for object detection. *BMVC*, 2021. 2
- [35] Gaoang Wang, Jenq-Neng Hwang, Craig Rose, and Farron Wallace. Uncertainty-based active learning via sparse modeling for image classification. *Transactions on Image Processing*, 2018. 3
- [36] Mengde Xu, Zheng Zhang, Han Hu, Jianfeng Wang, Lijuan Wang, Fangyun Wei, Xiang Bai, and Zicheng Liu. End-to-end semi-supervised object detection with soft teacher. In *ICCV*, 2021. 1
- [37] Donggeun Yoo and In So Kweon. Learning loss for active learning. In *CVPR*, 2019. 1, 3, 7
- [38] Weiping Yu, Sijie Zhu, Taojiannan Yang, and Chen Chen. Consistency-based active learning for object detection. In *CVPR*, 2022. 2, 7
- [39] Tianning Yuan, Fang Wan, Mengying Fu, Jianzhuang Liu, Songcen Xu, Xiangyang Ji, and Qixiang Ye. Multiple instance active learning for object detection. In *CVPR*, 2021. 2
- [40] Fedor Zhdanov. Diverse mini-batch active learning. *arXiv preprint arXiv:1901.05954*, 2019. 3, 7
- [41] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019. 6, 13

# Active Learning for Object Detection with Non-Redundant Informative Sampling

## Supplementary Material

### 1. Loss Inequality Eq. (2)

We denote the losses of a sample  $u$  before and after the training with a sample  $v$  as  $l(u; \theta_i)$  and  $l(u; \theta_{i+1})$ , and similarly for  $v$ . The  $\kappa$ -Lipschitz continuity of the loss function  $l(\cdot)$  (we refer to the Core-Set [29] for proof) provides the following inequality for any points  $u$  and  $v$ ,

$$|l(u; \theta) - l(v; \theta)| \leq \kappa * d(u, v) \quad (11)$$

Before the training of  $v$ , we have the inequality,

$$|l(u; \theta_i) - l(v; \theta_i)| \leq \kappa * d(u, v) \quad (12)$$

Since we select the sample with a higher loss, we assume that the initial loss at  $v$  is higher than the initial loss at  $u$ ,

$$l(v; \theta_i) - l(u; \theta_i) \leq \kappa * d(u, v) \quad (13)$$

Similarly, for the final weights  $\theta_{i+1}$ ,

$$|l(u; \theta_{i+1}) - l(v; \theta_{i+1})| \leq \kappa * d(u, v) \quad (14)$$

Following Core-Set [29], we assume the final loss at  $v$ ,  $l(v; \theta_{i+1})$ , to be minimized, approaching zero. Therefore, after the training of  $v$ , we have the inequality,

$$l(u; \theta_{i+1}) \leq \kappa * d(u, v) \quad (15)$$

Summing up the two inequalities in Eq. (13), Eq. (15),

$$l(v; \theta_i) - l(u; \theta_i) + l(u; \theta_{i+1}) \leq 2 * \kappa * d(u, v) \quad (16)$$

Isolating  $l(u; \theta_{i+1})$ ,

$$l(u; \theta_{i+1}) \leq l(u; \theta_i) + 2 * \kappa * d(u, v) - l(v; \theta_i) \quad (17)$$

### 2. Hyperparameter $\lambda$

As stated in the main manuscript, we use  $\lambda = \alpha \cdot d_{\max}$  for linear similarity and  $\lambda = \frac{(\alpha \cdot d_{\max})^2}{\pi}$  for Gaussian similarity. In this section, we explain the reasoning behind our choice of defining the Gaussian similarity parameter in this manner.

Let  $\lambda_G$  and  $\lambda_l$  denote the hyperparameter for the Gaussian and linear similarity function respectively. Our goal is to choose  $\lambda_G$  such that the Gaussian and the linear similarity function have a similar behavior. We argue that this happens when the integral between the two functions is 0:

$$\int_0^{\infty} e^{-\frac{1}{\lambda_G} s^2} - \max \left\{ 0, 1 - \frac{s}{\lambda_l} \right\} ds \stackrel{!}{=} 0 \quad (18)$$

This is equivalent to require:

$$\int_0^{\infty} e^{-\frac{1}{\lambda_G} s^2} ds \stackrel{!}{=} \int_0^{\infty} \max \left\{ 0, 1 - \frac{s}{\lambda_l} \right\} ds \quad (19)$$

Let's calculate both sides. We know for a fact that the integral over the real numbers of the probability density function of a Gaussian distribution is 1:

$$\frac{1}{\sigma_G \sqrt{2\pi}} \int_{\mathbb{R}} e^{-\frac{s^2}{2\sigma_G^2}} ds = 1 \text{ for all } \sigma_G > 0$$

By the change of variable  $\lambda_G = 2\sigma_G^2$ , we obtain

$$\int_{\mathbb{R}} e^{-\frac{s^2}{\lambda_G}} ds = \sqrt{\lambda_G \pi} \text{ and thus } \int_0^{\infty} e^{-\frac{s^2}{\lambda_G}} ds = \frac{\sqrt{\lambda_G \pi}}{2}$$

for the left hand side of Eq. (19) by the symmetry. Next, we simplify the right hand side:

$$\int_0^{\infty} \max \left\{ 0, 1 - \frac{s}{\lambda_l} \right\} ds = \int_0^{\lambda_l} 1 - \frac{s}{\lambda_l} ds = \left( s - \frac{s^2}{2\lambda_l} \right) \Big|_0^{\lambda_l} = \frac{\lambda_l}{2}$$

From the calculations, we infer that the integrals in Eq. (19) are equal when  $\frac{\sqrt{\lambda_G \pi}}{2} = \frac{\lambda_l}{2}$  which can be achieved by setting

$$\lambda_G = \frac{\lambda_l^2}{\pi}. \quad (20)$$

Intuitively, this implies that the Gaussian similarity score is  $e^{-\pi} \approx 0.043$  at a distance of  $\lambda_l$ . When using the Gaussian similarity, we therefore suggest fixing a value  $\alpha \in (0, 1]$  and scaling to  $\lambda = \frac{(\alpha \cdot d_{\max})^2}{\pi}$  in each AL cycle for hyperparameter tuning. An advantage of this theoretical investigation is that once a good hyperparameter choice  $\lambda_l > 0$  for the linear similarity function is found, we can act on the assumption that setting  $\lambda_G = \frac{\lambda_l^2}{\pi}$  in the Gaussian similarity function performs similar and vice versa.

**Hyperparameter optimization for  $\alpha$ .** As explained in the manuscript, we select a fixed value for  $\alpha$  and scale it with  $d_{\max}$  to derive the corresponding  $\lambda$  value. Thus, determining the optimal value of  $\alpha$  is critical in our work. We

Variant	Similarity	$\alpha$										
		0.01	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
NORIS-Max	linear	39.4	39.7	39.6	39.8	39.7	39.0	39.8	39.7	<b>40.1</b>	39.5	40.0
	gaussian	40.6	40.4	<b>40.8</b>	40.3	40.5	40.5	39.9	40.2	40.3	40.3	39.7
NORIS-Sum	linear	39.6	39.5	39.1	39.0	39.5	39.8	40.2	39.9	39.5	40.1	<b>40.4</b>
	gaussian	40.0	40.4	40.2	40.5	40.6	<b>41.0</b>	40.3	40.6	40.4	40.0	40.7

Table 3. Hyperparameter search for two variants of NORIS and the two similarity measures. Results are reported on the KITTI dataset with the mAP metric. Bold indicates the best-performing parameter value for the row.

present the results of our parameter search in Tab. 3. Notably, we observe that higher values of  $\alpha$  (0.8 and 1.0) tend to yield superior performance for the linear similarity. Conversely, for the Gaussian similarity, the most effective values are 0.2 and 0.5 for the NORIS-Max and NORIS-Sum variants, respectively.

### 3. NORIS-Max Algorithm

We present the NORIS-Max algorithm in Algorithm 2. In each iteration of the while-loop, the algorithm selects the sample  $u_{next}$  that leads to the highest immediate gain and adds it to  $S$ . Subsequently, for the most similar sample  $u_c = \arg \max_{u \in X_U \setminus S} (sim(u, u_{next}))$ , the information score is updated by subtracting the uncertainty of  $u_{next}$  scaled by the similarity  $sim(u, u_{next})$ . The runtime complexity of this algorithm is  $\mathcal{O}(B \cdot |X_U|)$ .

---

#### Algorithm 2 NORIS-Max

---

**Require:** unlabeled dataset  $X_U$ , batch size  $1 \leq B \leq |X_U|$

$S \leftarrow \emptyset$

**while**  $|S| \neq B$  **do**

$u_{next} \leftarrow \arg \max_{u \in X_U \setminus S} \sigma(u)$

$S \leftarrow S \cup \{u_{next}\}$

$u_c \leftarrow \arg \max_{u \in X_U \setminus S} (sim(u, u_{next}))$

$\sigma(u_c) \leftarrow \sigma(u_c) - sim(u_c, u_{next}) * \sigma(u_{next})$

**end while**

---

### 4. Extending MAL and VAAL to Object Features

We provide a comprehensive description of how we incorporate object features into the VAAL [32] and MAL [10] methods. Figure 6 demonstrates the process of extracting object features from the CenterNet detector. Following object detection, we crop the corresponding bounding box locations from the intermediate feature map to obtain object features, as described in the main manuscript.

Fig. 7a and Fig. 7b depict the utilization of object features for the MAL and VAAL methods, respectively. For the MAL approach, we input the object features into the

discriminator, and for the VAAL method, we feed the object features to a variational autoencoder (VAE). In both cases, we employ the discriminator predictions as a similarity score to the labeled set and rank images based on the most dissimilar object it contains.

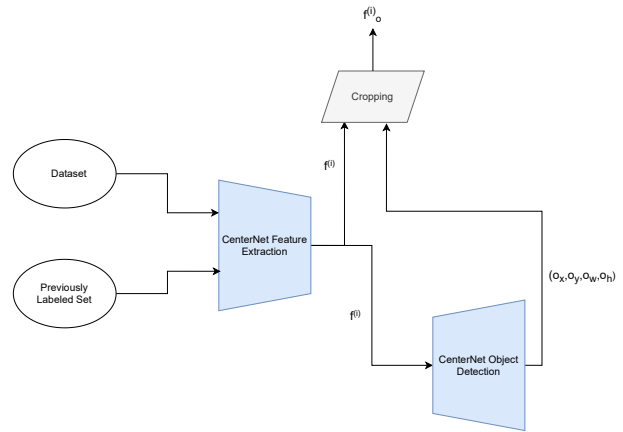


Figure 6. Illustration of extracting object features from the CenterNet detector.

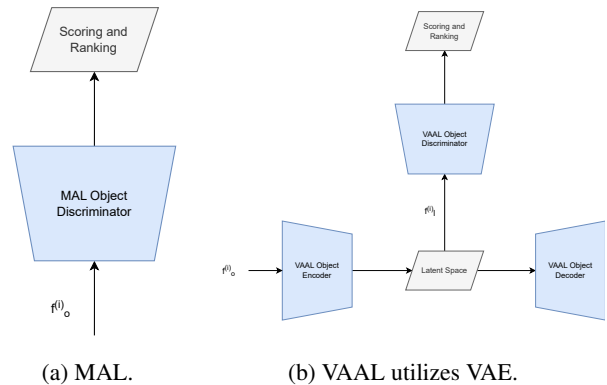


Figure 7. Illustration of the VAAL and MAL with object features.



## 5. CenterNet Features

In Fig. 2, we describe the process of extracting object features. An important consideration is determining which layer from the detector represents the intermediate features. In Fig. 8, we provide a detailed illustration of different intermediate features from the CenterNet detector [41]. We denote the features extracted from the detector as  $f_1$ ,  $f_2$ ,  $f_3$ ,  $f_4$ , with  $f_1$  being the closest feature to the input image and  $f_4$  being the feature closest to the output as depicted in Fig. 8.

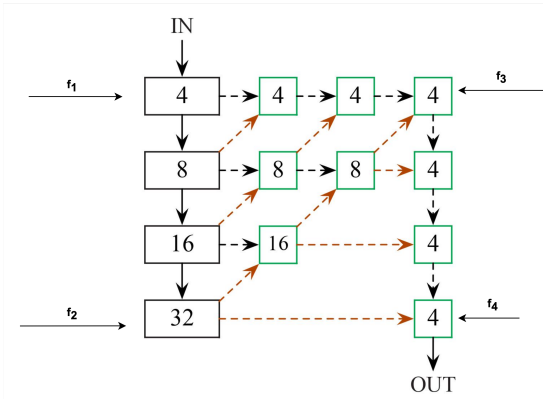


Figure 8. Illustration showcasing the intermediate network layers employed in our experiments. The object features are extracted from the layers represented in this visualization.

Then, in Fig. 9, we perform an ablation study to investigate the impact of using different intermediate feature maps from various layers of the CenterNet detector in our selection strategy. The results indicate that using features from the later layers leads to improved performance in our selection strategy. This finding indicates that utilizing task-specific features can be advantageous.

## 6. Ablation of the distance metric, aggregation function, and use of image features.

To investigate the impact of different design choices on the performance of our approach, we compare using different distance metrics, aggregation functions, and features. Specifically, we compare using cosine or Euclidean distance as the distance metric, aggregating object features into an image score using averaging or maximum, and using image features or not. Our results in Tab. 4 indicate that using image features leads to higher performance than not using them. The choice of the aggregation function has a negligible impact on the outcome, but we observe a slight improvement when using the maximum over averaging. Our results show a slight advantage of using Euclidean distance over cosine distance as the distance metric. We choose to use Euclidean distance, maximum aggregation and include image features for our experiments, as this configuration cor-

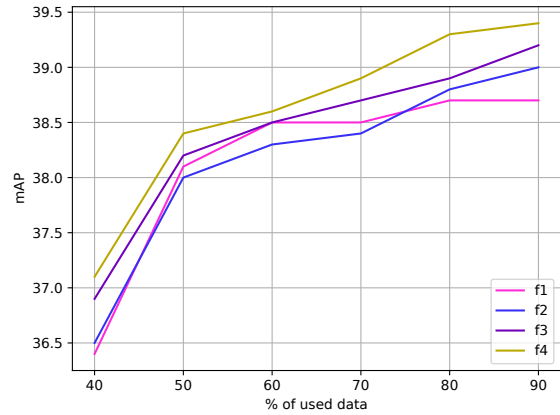


Figure 9. Comparison of using feature maps from different layers of the detector on KITTI.  $f_1$  represent the closest feature to the input and  $f_4$  represent the furthest.

responds to the last row with the highest performance in our results.

## 7. Exact Values from AL Figures and Variances

Due to space constraints in the main paper, we present our AL comparisons as plots. However, in the supplementary material, we provide the precise values and variances for our experiments. Tab. 5, Tab. 6, Tab. 7 and Tab. 8 provides the exact metric values for Figures 3a, 3b, 3c and 3d from the main paper. The mean and variances of three experiments trained with different random initializations are presented.

Dist.	Agg.	$f_I$	40%	60%	80%	100%
cos	avg		36.6	38.3	38.6	38.6
euc	avg		37.0	38.1	38.5	39.1
cos	max		36.7	38.1	38.6	39.0
euc	max		36.4	38.0	39.0	39.3
cos	avg	✓	36.9	37.9	38.7	39.4
euc	avg	✓	37.0	38.5	39.2	39.5
cos	max	✓	37.0	38.3	39.0	39.4
euc	max	✓	37.2	38.2	39.4	39.6

Table 4. Comparison of distance metric, aggregation function and using image features on KITTI.

	40	50	60	70	80	90
Random	60.80±0.28	64.10±0.27	66.10±0.27	68.10±0.18	68.80±0.37	69.20±0.23
CDAL	62.40±0.32	64.80±0.34	67.80±0.16	68.80±0.36	69.80±0.21	70.50±0.23
LL4AL	60.60±0.35	64.40±0.28	66.70±0.16	68.90±0.24	69.90±0.25	70.80±0.29
DBAL	61.60±0.32	65.20±0.22	67.70±0.19	69.50±0.37	71.10±0.29	71.50±0.29
CALD	62.90±0.16	66.00±0.34	67.80±0.18	69.10±0.33	70.40±0.35	71.20±0.16
NORIS	63.90±0.19	67.50±0.26	68.50±0.32	70.10±0.36	71.60±0.38	72.40±0.20

Table 5. Comparison of mAP with SOTA AL methods on PASCAL-VOC 2007. (Fig. 3a)

	40	50	60	70	80	90
Random	36.10±0.17	37.30±0.11	37.90±0.06	38.40±0.06	38.70±0.11	38.80±0.16
CALD	36.90±0.10	38.20±0.08	38.60±0.12	38.90±0.15	39.10±0.16	39.20±0.20
LL4AL	35.60±0.08	37.70±0.07	38.50±0.10	38.80±0.14	39.30±0.12	39.40±0.16
DBAL	36.00±0.12	37.80±0.17	38.60±0.14	39.60±0.06	39.80±0.13	40.00±0.06
CALD	36.70±0.10	38.50±0.05	39.20±0.16	39.40±0.11	39.50±0.10	39.60±0.14
NORIS	37.70±0.08	38.80±0.13	39.50±0.13	39.70±0.12	40.20±0.10	40.80±0.12

Table 6. Comparison of mAP with SOTA AL methods on KITTI *val.* (Fig. 3b)

	4.0	6.0	8.0	10.0	12.0	14.0	16.0
Random	59.30±0.54	63.40±0.33	67.00±0.75	70.60±0.56	71.70±0.58	73.00±0.70	74.20±0.69
CDAL	58.40±0.75	63.60±0.45	67.90±0.34	69.80±0.65	72.50±0.45	73.50±0.31	75.50±0.39
LL4AL	59.10±0.77	64.90±0.71	69.30±0.71	70.90±0.76	73.30±0.70	74.30±0.74	76.00±0.44
DBAL	59.30±0.67	65.20±0.76	68.40±0.32	73.10±0.69	74.30±0.32	76.70±0.35	77.30±0.45
LC	58.90±0.31	64.80±0.48	68.80±0.65	71.50±0.33	73.30±0.65	74.90±0.39	77.50±0.73
NORIS	61.20±0.55	66.80±0.69	70.60±0.71	74.20±0.66	75.80±0.69	76.90±0.44	78.80±0.49

Table 7. Comparison of mAP with SOTA AL methods on CIFAR-10. (Fig. 3c)

	4.0	6.0	8.0	10.0	12.0	14.0	16.0
Random	41.90±0.33	43.10±0.17	45.70±0.13	46.90±0.37	47.90±0.31	48.40±0.24	49.40±0.10
CDAL	41.60±0.16	43.50±0.28	45.80±0.26	47.50±0.14	48.80±0.38	49.40±0.12	49.90±0.34
LL4AL	41.60±0.24	43.90±0.27	46.50±0.37	48.00±0.10	48.90±0.14	49.00±0.25	49.80±0.40
DBAL	41.10±0.10	44.90±0.33	47.20±0.25	47.90±0.38	49.10±0.14	50.10±0.10	50.80±0.14
LC	41.70±0.19	44.40±0.10	45.80±0.18	47.80±0.12	48.60±0.20	49.40±0.23	50.00±0.10
NORIS	41.90±0.15	45.50±0.25	47.70±0.38	48.30±0.40	49.50±0.13	50.70±0.32	51.00±0.16

Table 8. Comparison of mAP with SOTA AL methods on CIFAR-100. (Fig. 3d)