

---

# Inconsistency-Based Data-Centric Active Open-Set Annotation

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Active learning is a commonly used approach that reduces the labeling effort  
2 required for training deep neural networks. However, the effectiveness of current  
3 active learning methods is limited by their closed-world assumptions, which assume  
4 that all data in the unlabelled pool comes from a set of predefined known classes.  
5 This assumption is often not valid in practical situations, as there may be unknown  
6 classes in the unlabeled data, leading to the active open-set annotation problem. The  
7 presence of unknown classes in the data can significantly impact the performance  
8 of existing active learning methods due to the uncertainty they introduce. To  
9 address this issue, we propose a novel data-centric active learning method called  
10 NEAT that actively annotates open-set data. NEAT is designed to label known  
11 classes data from a pool of both known and unknown classes unlabeled data. It  
12 utilizes the clusterability of labels to identify the known classes from the unlabeled  
13 pool and selects informative samples from those classes based on a consistency  
14 criterion that measures inconsistencies between model predictions and local feature  
15 distribution. Unlike the recently proposed learning-centric method for the same  
16 problem, NEAT is much more computationally efficient and is a data-centric active  
17 open-set annotation method. Our experiments demonstrate that NEAT achieves  
18 significantly better performance than state-of-the-art active learning methods for  
19 active open-set annotation.

20 

## 1 Introduction

21 The remarkable performance of modern deep neural networks owes much to the availability of  
22 large-scale datasets such as ImageNet [1]. However, creating such datasets is a challenging task that  
23 requires a significant amount of effort to annotate data points [2, 3, 4, 5]. Fortunately, active learning  
24 offers a solution by enabling us to label only the most *significant* samples [6, 7, 8, 9, 10]. In active  
25 learning, a small set of labeled samples from known classes is combined with a pool of unlabeled  
26 samples, and the objective is to identify which samples to label from the pool for better training the  
27 model. Typical active learning methods for training deep neural networks involve selecting samples  
28 that have high levels of uncertainty [6, 11, 12], are in close proximity to the classification boundary  
29 [13], or unveil cluster structures from the data [14, 15].

30 Despite the considerable body of research on active learning, its practical implementation in an *open-*  
31 *world* context remains relatively unexplored [16]. Unlike in closed-world active learning settings,  
32 where the unlabeled data pool is assumed to consist only of known classes, the open world introduces  
33 unknown classes into the unlabeled data pool, making active open-set annotations a challenge. For  
34 instance, when collecting data to train a classifier to distinguish between different dog breeds as  
35 known classes, a large number of unlabeled images collected from online sources may include images  
36 from unknown classes such as wolves and coyotes. The task is to identify images with known classes  
37 and select informative samples for labeling, while avoiding samples from unknown classes. Existing

38 active learning methods face a significant obstacle in this task, as they may tend to select samples  
39 from unknown classes for labeling due to their high uncertainty [16], which is undesirable.

40 Addressing the active open-set annotation problem requires specialized methods, and one such  
41 approach called LFOSA is proposed recently by [16]. LFOSA is a learning-based active open-set  
42 annotation as it involves training a detector network with an additional output for unknown classes.  
43 The predictions of the detector network on unlabeled data are used for identifying known classes.  
44 While this approach can achieve impressive performance, it has two limitations: 1) training the  
45 additional detector network is costly, 2) and it is difficult to identify informative samples from the  
46 known classes as there is a contradiction in that while excluding unknown classes is necessary, it is  
47 also easy to exclude informative samples from the known classes.

48 In this paper, we propose a novel inconsistency-based data-centric active learning method to actively  
49 annotate informative samples in an open world, which not only reduces computation cost but also  
50 improves the performance of active open-set annotation. Rather than using a learning-based approach  
51 to differentiate known and unknown classes, we suggest a data-centric perspective that naturally  
52 separates them by label clusterability, eliminating the need for an additional detector network. In  
53 addition, our method involves selecting informative samples from known classes by estimating the  
54 inconsistency between the model’s prediction and local feature distribution. For example, suppose the  
55 model predicts an unlabeled sample as a wolf, but the majority of nearby samples are actually dogs. In  
56 that case, we would choose to label this unlabeled sample. The proposed inconsistency-based active  
57 learning approach shares a similar spirit to the version-space-based approach [11, 12, 17]. However,  
58 there are two key differences. Firstly, our hypothesis class consists of deep neural networks. Secondly,  
59 the version-space-based approach identifies uncertain examples by analyzing the consistencies among  
60 multiple models, our approach leverages a fixed model and estimates the consistencies between the  
61 model prediction and the local feature distribution.

62 In summary, the contributions are the paper are as follows,

- 63 • We propose a novel and efficient inconsistency-based data-centric active learning method,  
64 called NEAT, for selecting informative known classes samples from a pool of both known  
65 and unknown classes.
- 66 • Compared with the learning-based method, the proposed data-centric active learning method  
67 is computationally efficient and can effectively identify informative samples from the known  
68 classes.
- 69 • Extensive experiments show that NEAT achieves much better results compared with the  
70 standard active learning methods and the method which is specifically designed for active  
71 open-set annotation. In particular, NEAT achieves an average accuracy improvement of  
72 9% on **CIFAR10**, **CIFAR100** and **Tiny-ImageNet** compared with existing active open-set  
73 annotation method with the same labelling budget.

## 74 2 Related Work

75 Active Learning, extensively studied in machine learning [18, 6, 7, 8, 9], focuses on selecting the most  
76 uncertain samples for labeling. Uncertainty-based active learning methods [18] employ measures  
77 such as entropy [19], least confident [6], and margin sampling [20]. Another widely used approach  
78 is the version-space-based method [11, 12, 17, 21], where multiple models consistent with labeled  
79 samples are maintained, and inconsistent predictions prompt label querying.

80 For deep neural networks, specific active learning methods have been developed. For instance,  
81 Core-Set [14] advocates for an active learning method that selects samples which are representative  
82 of the whole data distribution. Badge [15] selects samples based on predictive uncertainty and sample  
83 diversity. In particular, Badge leverages  $k$ -MEANS++ to select a set of samples which have diverse  
84 gradient magnitudes. The approach BGADL [22] integrates active learning and data augmentation.  
85 It leverages Bayesian inference in the generative model to create new training samples from selected  
86 existing samples. These generated samples are then utilized to enhance the classification accuracy  
87 of deep neural networks. CEAL [23] combines pseudo-labeling of clearly classified samples and  
88 actively labeled informative samples to train deep neural networks. For deep object detection, active  
89 learning is introduced in [24], utilizing prediction margin to identify valuable instances for labeling.  
90 DFAL [13] is an adversarial active learning method for deep neural networks that selects samples

91 based on their distance to the decision boundary, approximated using adversarial examples. The  
 92 samples closest to the boundary are chosen for classifier training. Recently, LFOSA [16] is proposed  
 93 as the first active learning method for active open-set annotation. LFOSA trains a detector network  
 94 to identify known classes and selects samples based on model confidence.

### 95 3 Preliminary

#### 96 3.1 Active Open-Set Annotation

97 Assume the sample-label pair  $(X, Y)$  has the joint distribution  $P_{XY}$ . For a given sample  $\mathbf{x}$ , the label  
 98  $y$  can be determined via the conditional expectation,

$$\eta(\mathbf{x}) = \mathbb{E}[Y|X = \mathbf{x}] \quad (1)$$

99 We consider a pool-based active learning setup. We begin by randomly sampling a small labeled  
 100 set denoted as  $L = \{\mathbf{x}_i, y_i\}_{i=1}^N$  from the joint distribution  $P_{XY}$ , where  $\mathbf{x}_i \in X$  and  $y_i \in Y$ . We  
 101 also have access to a large set of unlabeled samples, denoted as  $U$ . Given a deep neural network  
 102  $f_\theta$  parameterized by  $\theta$ , the expected risk of the network is computed as  $R = \mathbb{E}[\ell(y, f_\theta(\mathbf{x}))] =$   
 103  $\int l(y, f_\theta(\mathbf{x}))dP_{XY}$ , where  $\ell(y, f_\theta(\mathbf{x}))$  represents the classification error. We conduct  $T$  query rounds,  
 104 where in each round  $t$ , we are allotted a fixed labeling budget  $B$  to identify  $B$  informative samples  
 105 from  $U$ , denoted as  $U_B$ . The queried samples  $U_B$  are given labels and added to the initial labeled  
 106 set, with the aim of minimizing the sum of cross-entropy losses  $\sum_{(\mathbf{x}, y) \in L \cup U_B} \ell_{CE}(f_\theta(\mathbf{x}), y)$  and  
 107 reducing the expected risk of the model  $f_\theta$  [6]. The term  $\ell_{CE}(f_\theta(\mathbf{x}), y)$  represents the cross-entropy  
 108 loss between the predicted and true labels for a sample  $(\mathbf{x}, y)$  [25].

109 In the context of active open-set annotation [16], the unlabeled samples in  $U$  include both data from  
 110 known classes  $Y$  and unknown classes  $\tilde{Y}$ . In particular,  $Y \cap \tilde{Y} = \emptyset$ . Active open-set annotation  
 111 presents a greater challenge than standard active learning, as it requires the active learning method  
 112 1) to distinguish between known and unknown classes from the unlabeled samples 2) and to select  
 113 informative samples exclusively from the known classes.

#### 114 3.2 Learning-Based Active Open-Set Annotation

115 The study of active open-set annotation is currently limited. A recent paper [16] proposed a new  
 116 learning-based method called LFOSA which is specially designed for active open-set annotation.  
 117 This approach combines a model,  $f_\theta$ , trained on known classes with a detector network that identifies  
 118 unknown class samples. Suppose there are  $C$  known classes, the detector network has  $C + 1$  outputs,  
 119 similar to OPENMAX[26], which allow for the classification of unlabeled data into known and  
 120 unknown classes. During an active query, LFOSA focuses only on unlabeled data classified as known  
 121 classes, and then uses a Gaussian Mixture Model to cluster the activation values of detected known  
 122 classes data into two clusters. The data closer to the larger cluster means is selected for labeling. This  
 123 learning-based approach has shown significant improvement over traditional active learning methods  
 124 for active open-set annotation. However, there are still challenges that need to be addressed, such as  
 125 the additional computation cost required to train the detector network and the difficulty of identifying  
 126 informative samples from the known classes.

### 127 4 NEAT

128 Algorithm 1 describes the detailed algorithm of NEAT. Initially, NEAT begins with a randomly drawn  
 129 labeled set  $L$  from known classes. In each query round, NEAT performs two main steps. Firstly, it  
 130 identifies unlabeled samples whose neighbors in the labeled set belong to known classes. Secondly,  
 131 NEAT selects a batch of unlabeled samples for labeling by estimating the inconsistency between the  
 132 model prediction and local feature distribution. Unlike the learning-based method discussed in [16],  
 133 NEAT is computationally efficient. Moreover, it effectively decouples the detection of known classes  
 134 from the identification of informative samples, enabling it to identify informative samples even from  
 135 the known classes, which is a challenge for existing active learning methods.

---

**Algorithm 1** NEAT: Inconsistency-Based Data-Centric Active Open-Set annotation.

---

**Require:** A deep neural network  $f_\theta$ , initial labeled set  $L$ , a set of known class  $Y_{\text{Known}}$ , unlabeled labeled set  $U$ , number of query rounds  $T$ , number of examples in each query batch  $B$ , a pre-trained model  $M$ , number of neighbors  $K$ .

- 1: Use the pre-trained model  $M$  for extracting features on  $L$  and  $U$ .
- 2: **for**  $t \leftarrow 1$  to  $T$  **do**
- 3:     Train the model  $f_\theta$  on  $L$  by minimizing  $\sum_{(x,y) \in L} \ell_{\text{CE}}(f_\theta(x), y)$
- 4:      $S \leftarrow \{\}$
- 5:     For each sample  $x \in U$ :
- 6:         Compute the output of the softmax function as  $P_x$
- 7:         Find the  $K$ -nearest neighbors  $\{N_1(x), N_2(x), \dots, N_K(x)\}$  of  $x$  in  $L$  based on the extracted features
- 8:         If all the labels of  $\{N_1(x), N_2(x), \dots, N_K(x)\}$  belong to known classes  $Y_{\text{Known}}$ , then  
         $S \leftarrow S \cup \{x\}$
- 9:         Compute the score  $I(x)$  using Eq. 3 for each sample  $x \in S$
- 10:         Rank the samples based on  $I(x)$  and denote the  $B$  samples which have the largest scores as  $U_B$ .
- 11:         Query the labels of each sample in  $U_B$ .
- 12:          $U \leftarrow U \setminus U_B, L \leftarrow L \cup U_B$
- 13: **end for**

---

136 **4.1 Data-Centric Known Class Detection**

137 In the learning-based method [16], an additional detector network is trained to differentiate known  
138 classes and unknown classes. However, the training cost is high. Instead, NEAT takes a data-centric  
139 perspective for finding known classes samples based on feature similarity. In particular, NEAT relies  
140 on label clusterability to find known classes from the unlabeled pool.

141 **Label clusterability.** We leverage the intuition that samples with similar features should belong to  
142 the same class [27, 28, 29]. This intuition can be formally defined as follows,

143 **Definition 4.1.**  $((K, \sigma_K)$  label clusterability). A dataset  $D$  satisfies  $(K, \sigma_K)$  label clusterability if  
144 for all  $x \in D$ , the sample  $x$  and its  $K$ -Nearest-Neighbors ( $K$ -NN)  $\{x_1, x_2, \dots, x_K\}$  belong to the  
145 same label with probability at least  $1 - \sigma_K$ .

146 When  $\sigma_K = 0$ , then it is called  $K$ -NN label clusterability [28]. Recent methods for noisy label  
147 learning [28, 29] leverage label clusterability for detecting examples with noisy labels. The label  
148 clusterability of the dataset is closely related to the smoothness condition introduced in [30] for  
149 analyzing nearest neighbor classifier [30] and non-parametric active learning [31],

150 **Definition 4.2.**  $((\alpha, L)$ -smooth) [30]. The conditional expectation  $\eta$  is  $(\alpha, L)$ -smooth if for all  $x$  and  
151  $x' \in D$ ,

$$|\eta(x) - \eta(x')| \leq L\mu(B(x, \rho(x, x')))^\alpha \quad (2)$$

152 where  $B(x, \rho(x, x'))$  is an open-ball centred at  $x$  with radius  $\rho(x, x')$ ,  $\mu$  is a measure on the input  
153 space.

154 When a conditional expectation is  $(\alpha, L)$ -smooth, it intuitively implies that similar samples are highly  
155 likely to have the same label. Therefore, if the features accurately capture the semantic similarity  
156 between samples, the similarity in features reflects the clusterability of labels.

157 **Feature extraction.** To utilize the clusterability of labels for identifying known classes, we need  
158 to extract features from unlabelled inputs that group semantically similar samples in the feature  
159 space. Additionally, the quality of these features will inevitably impact the clusterability of the labels.  
160 Instead of developing a separate detector, as demonstrated in [16], we suggest taking advantage of  
161 pre-trained large language models for feature extraction, which have been demonstrated to possess  
162 exceptional zero-shot learning ability [32]. In particular, we leverage CLIP [32] to extract features for  
163 both the labeled and unlabeled data, providing high-quality features for calculating feature similarity.

164 **Known class detection.** By utilizing the features extracted by CLIP from both the labeled set  $L$   
165 and the unlabeled set  $U$ , we can identify the  $K$ -nearest neighbors  $\{N_1(x), N_2(x), \dots, N_K(x)\}$  in the

166 labeled set  $L$  for each sample  $\mathbf{x} \in U$ , using cosine distance. Each  $N_k(\mathbf{x}) \in L$  represents the  $k$ -th  
 167 closest samples in  $L$  to the unlabeled sample  $\mathbf{x}$ . Afterward, we compute the count of neighbors with  
 168 known and unknown classes for each unlabeled sample  $\mathbf{x}$ , assuming label clusterability. If there are  
 169 many unknown-class samples close to an unlabeled sample, it is more likely to belong to a known  
 170 class. Finally, we determine the unlabeled samples with all neighbors belonging to known classes as  
 171 potential known-class samples.

## 172 4.2 Inconsistency-Based Active Learning

173 To improve accuracy with a fixed labeling budget, it is crucial to select informative samples for  
 174 labeling. One active learning strategy, motivated by theory, is the version-space-based approach  
 175 [11, 12, 17, 21]. This approach involves maintaining a set of models that are consistent with the  
 176 current labeled data, and an unlabeled sample is selected for labeling if two models produce different  
 177 predictions. However, implementing this approach for deep neural networks is challenging due  
 178 to the computational cost of training multiple models [15]. To address this issue, we propose  
 179 an inconsistency-based active learning method that does not require training multiple models and  
 180 naturally leverages features produced by CLIP.

181 Given an unlabeled sample  $\mathbf{x} \in U$ , the model  $f_\theta$ 's prediction for  $\mathbf{x}$ , denoted as  $P_{\mathbf{x}} \in \mathbb{R}^C$ , represents  
 182 a probability vector where  $P_{\mathbf{x}}[c]$  is the model's confidence that  $\mathbf{x}$  belongs to class  $c$ . Since we  
 183 lack ground-truth labels, measuring prediction accuracy is impossible. To address this, we propose  
 184 evaluating the importance of the sample for improving model training by assessing whether the model  
 185 prediction aligns with local feature similarity. For example, if the model predicts an unlabeled sample  
 186 as a dog with low probability but the majority of the sample's neighbors belong to the dog class, then  
 187 either the model's prediction is incorrect or the sample is near the decision boundary between the dog  
 188 class and the true class. In either case, the sample can be labeled to improve model training.

189 Given the  $K$ -nearest neighbors  $\{N_1(\mathbf{x}), N_2(\mathbf{x}), \dots, N_K(\mathbf{x})\}$  of the example  $\mathbf{x}$ , we first construct a  
 190 vector  $V_{\mathbf{x}} \in \mathbb{R}^C$  with  $V_{\mathbf{x}}[c] = \sum_k \mathbf{1}(Y_k(\mathbf{x}) = c)$ , where  $Y_k(\mathbf{x})$  is the label of  $k$ -th nearest neighbor  
 191 of  $\mathbf{x}$  and  $\mathbf{1}(\cdot)$  is an indicator function. Then the vector  $V_{\mathbf{x}}$  is normalized via the softmax function  
 192 to be a probabilistic vector  $\tilde{V}_{\mathbf{x}}$ . The inconsistency between the model prediction and local feature  
 193 similarity is computed using cross-entropy as,

$$I(\mathbf{x}) = - \sum_{c=1}^C P_{\mathbf{x}}[c] \log \tilde{V}_{\mathbf{x}}[c]. \quad (3)$$

194 A large  $I(\mathbf{x})$  indicates that the model prediction is inconsistent with local feature similarity, similar  
 195 to version-space-based approach, the unlabeled sample is selected for labeling. In each query round,  
 196 we rank all the identified known classes samples in the first stage using  $I(\mathbf{x})$  and select the top  $B$   
 197 samples for labeling.

## 198 5 Experiments

### 199 5.1 Experimental Settings and Evaluation Protocol

200 **Datasets and models.** We consider **CIFAR10** [33], **CIFAR100** [33], and **Tiny-Imagenet** [34], to  
 201 evaluate the performance of our proposed method. Similar to existing methods for active open-set  
 202 annotations [16], we leverage a ResNet-18 [35] architecture to train the classifier for the known  
 203 classes. For the proposed method, we leverage CLIP [32] to extract the features for both the known  
 204 classes and unknown classes.

205 **Active open-set annotation.** In accordance with [16], the experiment randomly selects 40 classes,  
 206 20 classes, and 2 classes from **Tiny-Imagenet**, **CIFAR100**, and **CIFAR10**, respectively, to be known  
 207 classes, while the remaining classes are treated as unknown classes. To begin with, following [16],  
 208 8% of the data from the known classes in **Tiny-ImageNet** and **CIFAR100**, and 1% of the data from  
 209 the known classes in **CIFAR10** are randomly selected to form an initial labeled set. The rest of the  
 210 known class data and the unknown class data are combined to form the unlabeled data pool.

211 **Baseline methods.** We consider the following active learning methods as baselines,

- 212 1. RANDOM: The naive baseline which randomly selects samples for annotation;

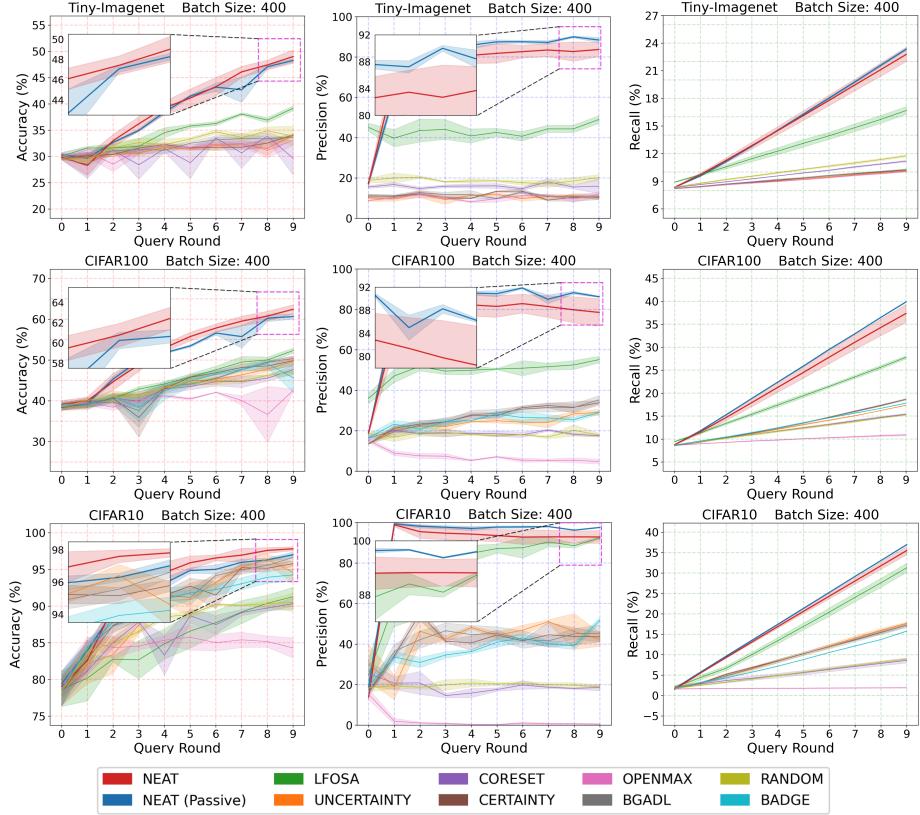


Figure 1: NEAT achieves higher precision, recall and accuracy compared with existing active learning methods for active open-set annotation. We evaluate NEAT and the baseline active learning methods on **CIFAR10**, **CIFAR100** and **Tiny-ImageNet** based on accuracy, precision and recall.

- 213     2. UNCERTAINTY [6]: A commonly used active learning method which selects samples with  
214       the highest degree of uncertainty as measured by entropy;
- 215     3. CERTAINTY [6]: A common baseline for active learning which selects samples with the  
216       highest degree of certainty as measured by entropy;
- 217     4. CORESET [14]: This method identifies a compact, representative subset of training data for  
218       annotation;
- 219     5. BGADL [22]: A Bayesian active learning method which leverages generative to select  
220       informative samples;
- 221     6. OPENMAX [26]: A representative open-set classification method which can differentiate  
222       between known classes and unknown classes;
- 223     7. BADGE [15]: A active learning method designed for deep neural networks which select a  
224       batch of samples with diverse gradient magnitudes;
- 225     8. LFOUSA [16]: A learning-based active open-set annotation method which selects samples  
226       based on maximum activation value (MAV) modeled by a Gaussian Mixture Model;
- 227     9. NEAT (Passive): We also consider a baseline which is the passive version of NEAT, that is,  
228       we do not leverage the proposed inconsistency-based active learning methods for selecting  
229       samples for labeling and just randomly sample from the identified known classes samples.

230 **Metrics.** We evaluate various active learning methods based on three key metrics: accuracy, precision,  
231 and recall. Accuracy reflects how accurately the model makes predictions. To measure recall, we use  
232 the ratio of selected known class samples to the total number of known class samples present in the  
233 unlabeled pool, denoted as  $N_{\text{known}}^t$  and  $N_{\text{known}}^{\text{total}}$ , respectively, in the query round  $t$ . Precision, on the  
234 other hand, measures the proportion of true known class samples among the selected samples in each

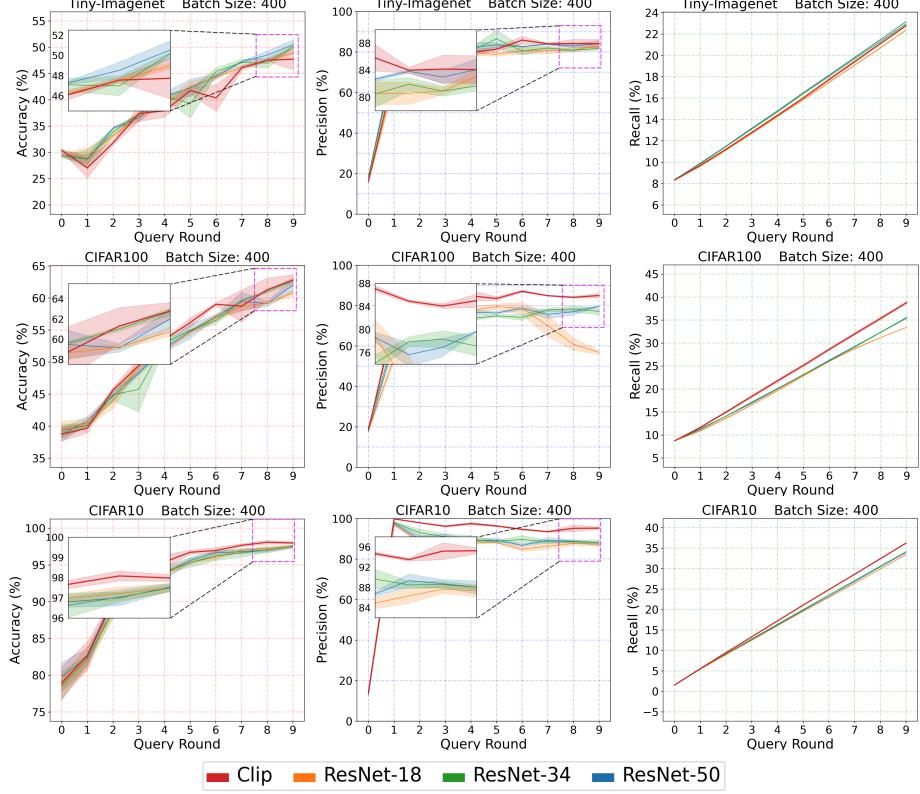


Figure 2: NEAT is robust to the choices of pre-trained models. We show the results of using different pre-trained models for feature extraction.

235 query round,

$$\text{precision} = \frac{N_{\text{known}}^t}{B} \quad \text{recall} = \frac{\sum_{t=1}^T N_{\text{known}}^t}{N_{\text{total known}}} \quad (4)$$

236 **Implementations details.** The classification model was trained for 100 epochs using SGD with  
237 an initial learning rate of 0.01. The learning rate decayed by 0.5 every 20 epochs, and the training  
238 batch size was set to 128. In total, there were 9 query rounds with a query batch size of 400. To  
239 ensure robustness, all experiments were repeated three times with different random seeds, and both  
240 the average results and standard deviations were reported. The experiments were conducted on four  
241 A5000 NVIDIA GPUs.

## 242 5.2 Results

### 243 5.2.1 NEAT VS. Baselines

244 We assess the performance of the proposed NEAT and compared methods by plotting curves with the  
245 number of queries increasing (Figure 1). It is evident that regardless of the datasets, NEAT consistently  
246 surpasses other methods in all cases. In particular, NEAT achieves much higher selection recall and  
247 precision compared with existing active learning methods which demonstrates the effectiveness of the  
248 data-centric known class detection. **1)** In terms of recall, NEAT consistently outperforms other active  
249 learning methods by a significant margin. Notably, on **CIFAR10**, **CIFAR100**, and **Tiny-ImageNet**,  
250 NEAT achieves improvements of 4%, 12%, and 7%, respectively, compared to LFOSA [16], which  
251 is a learning-based active open-set annotation method. **2)** In terms of precision, NEAT consistently  
252 maintains a higher selection precision than other baselines, with a noticeable gap. Importantly,  
253 NEAT’s ability to differentiate between known and unknown classes is significantly improved by  
254 adding examples from unknown classes in the first query round. **3)** In terms of accuracy, NEAT  
255 consistently outperforms LFOSA [16] across all datasets, including **CIFAR10**, **CIFAR100**, and  
256 **Tiny-ImageNet**, with improvements of 6%, 11%, and 10%, respectively. These results indicate that

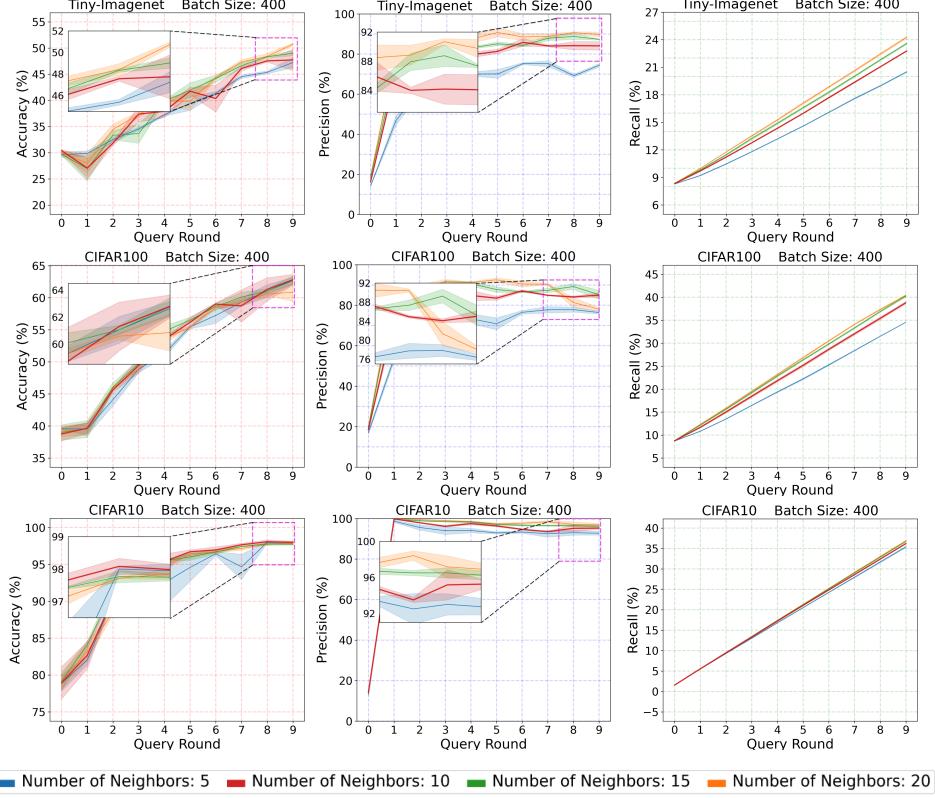


Figure 3: NEAT exhibits robustness to the choice of the number of neighbors, although a smaller value of  $K$  generally results in slightly lower performance.

the proposed NEAT method effectively addresses the open-set annotation (OSA) problem. It is worth noting that NEAT (Passive) achieves slightly higher recall and precision than NEAT. The possible reason is that by selecting informative samples from known classes, it is also more likely to include samples from unknown classes with high uncertainty. However, despite a slightly lower recall, NEAT achieves higher accuracy across all datasets compared to NEAT (Passive). This demonstrates that NEAT can effectively identify informative samples from known classes to train the model.

### 5.2.2 Ablation Studies

**Impact of feature quality.** We investigate the impact of different pre-trained models for feature extraction. The quality of the features is a crucial factor that may affect the label clusterability of the dataset. Specifically, we consider pre-trained ResNet-18, ResNet-34, and ResNet-50 on ImageNet, in addition to CLIP. Figure 2 demonstrates that NEAT achieves significantly better accuracy than the baseline active learning methods, regardless of the pre-trained models used. This highlights the robustness of NEAT in detecting known classes.

It's worth noting that while CLIP features achieve better accuracy on **CIFAR-10** and **CIFAR-100** than ResNets, the accuracy on **Tiny-ImageNet** is lower. This could be because the ResNets are pre-trained on ImageNet, and their features are better suited for Tiny-ImageNet. However, large language models like CLIP can generally provide high-quality features that are useful across different datasets.

**Influence of number of neighbors.** We further investigate how the number of neighbors impacts the detection of data-centric known classes. We consider different values of  $K$  (5, 10, 15, 20) and present the results in Figure 3. We observe that although a smaller value of  $K$  (e.g.,  $K = 5$ ) leads to slightly worse performance, other choices of  $K$  yield similar results. This suggests that a smaller  $K$  only captures the local feature distribution of the target sample, which may not provide a good characterization of the underlying feature space.

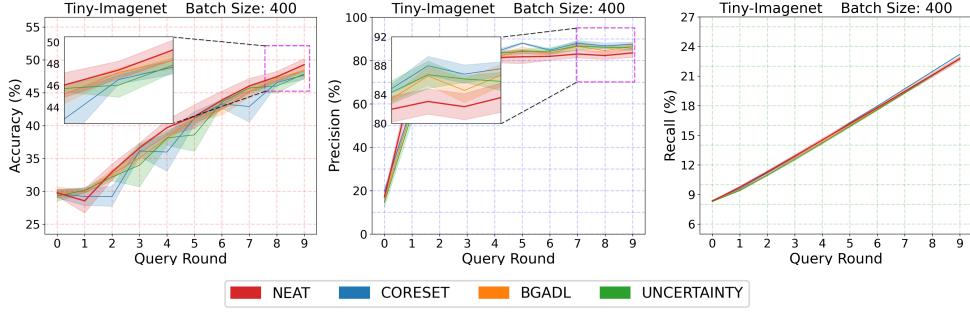


Figure 5: NEAT is effective compared with other active learning methods for deep neural networks.

281 **Effectiveness of inconsistency-based active learning.** We conducted an ablation study to evaluate  
 282 the effectiveness of the inconsistency-based active learning step in comparison to other active learning  
 283 methods for deep neural networks. We adopt a hybrid approach that combines *data-centric known*  
 284 *class detection* to initially identify potential known classes, followed by the application of different  
 285 active learning methods using the data from these detected known classes. The results on **Tiny-**  
 286 **ImageNet** are presented in Figure 5. The findings demonstrate that our proposed inconsistency-based  
 287 active learning method achieves higher accuracy when compared to other active learning methods.  
 288 These results indicate that the proposed inconsistency-  
 289 based active learning approach is capable of select-  
 290 ing more informative samples compared to alternative  
 291 methods.

292 Ablation studies of query batch size and the impact  
 293 of different classification models are included in the  
 294 Appendix.

295 **Visualization.** We used t-SNE [36] to visualize the  
 296 features produced by CLIP on CIFAR10, focusing on  
 297 the samples selected by NEAT and LFOSA [16]. We  
 298 randomly select a query round for plotting. Our results  
 299 show that LFOSA selects a large number of samples  
 300 from unknown classes, which explains its low precision.  
 301 On the other hand, almost all of the samples selected by  
 302 NEAT belong to the known classes, demonstrating its ef-  
 303 fectiveness in addressing the active open-set annotation  
 304 problem.

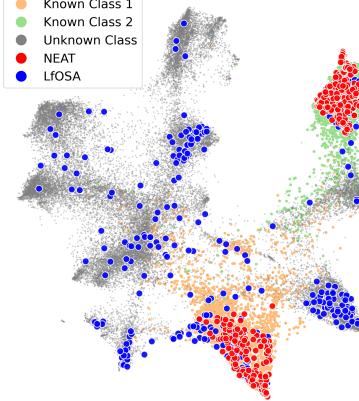


Figure 4: NEAT can accurately identify known classes from the unlabeled pool.

## 305 6 Summary

306 In this paper, we present a solution to the practical challenge of maintaining high recall when  
 307 identifying examples of known classes for target model training from a massive unlabeled open-set.  
 308 To address this challenge, we introduce a data-centric active learning method called NEAT, which  
 309 utilizes existing large language models for identifying known classes from the unlabeled pool. Our  
 310 proposed method offers several advantages over traditional active learning methods. Specifically,  
 311 NEAT uses the CLIP model to extract features, which eliminates the need for training a separate  
 312 detector. Additionally, our approach achieves improved results with low query numbers, resulting in  
 313 labeling cost savings. Furthermore, NEAT offers a plug-and-play solution with high adaptability to  
 314 various datasets.

315 NEAT has the potential to create a positive societal impact by reducing the labeling cost for learning  
 316 in an open world, which is crucial for practical applications such as autonomous driving and medical  
 317 imaging. We believe the proposed NEAT should not raise any ethical considerations. However, a  
 318 limitation of the proposed method is that for real-world application, maintaining pre-trained models  
 319 like CLIP can result in significant memory costs. Therefore, it may be necessary to quantize and  
 320 compress the model for more efficient deployment.

321 **References**

- 322 [1] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-  
323 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern  
324 recognition*, pages 248–255. Ieee, 2009.
- 325 [2] Alexander Sorokin and David Forsyth. Utility data annotation with amazon mechanical turk. In  
326 *2008 IEEE computer society conference on computer vision and pattern recognition workshops*,  
327 pages 1–8. IEEE, 2008.
- 328 [3] Rion Snow, Brendan O’connor, Dan Jurafsky, and Andrew Y Ng. Cheap and fast—but is it  
329 good? evaluating non-expert annotations for natural language tasks. In *Proceedings of the 2008  
330 conference on empirical methods in natural language processing*, pages 254–263, 2008.
- 331 [4] Vikas C Raykar, Shipeng Yu, Linda H Zhao, Gerardo Hermosillo Valadez, Charles Florin, Luca  
332 Bogoni, and Linda Moy. Learning from crowds. *Journal of machine learning research*, 11(4),  
333 2010.
- 334 [5] Peter Welinder, Steve Branson, Pietro Perona, and Serge Belongie. The multidimensional  
335 wisdom of crowds. *Advances in neural information processing systems*, 23, 2010.
- 336 [6] Burr Settles. Active learning literature survey. 2009.
- 337 [7] David A Cohn, Zoubin Ghahramani, and Michael I Jordan. Active learning with statistical  
338 models. *Journal of artificial intelligence research*, 4:129–145, 1996.
- 339 [8] Burr Settles. From theories to queries: Active learning in practice. In *Active learning and  
340 experimental design workshop in conjunction with AISTATS 2010*, pages 1–18. JMLR Workshop  
341 and Conference Proceedings, 2011.
- 342 [9] Alina Beygelzimer, Sanjoy Dasgupta, and John Langford. Importance weighted active learning.  
343 In *Proceedings of the 26th annual international conference on machine learning*, pages 49–56,  
344 2009.
- 345 [10] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image  
346 data. In *International conference on machine learning*, pages 1183–1192. PMLR, 2017.
- 347 [11] David Cohn, Les Atlas, and Richard Ladner. Improving generalization with active learning.  
348 *Machine learning*, 15:201–221, 1994.
- 349 [12] Maria-Florina Balcan, Alina Beygelzimer, and John Langford. Agnostic active learning. In  
350 *Proceedings of the 23rd international conference on Machine learning*, pages 65–72, 2006.
- 351 [13] Melanie Ducoffe and Frederic Precioso. Adversarial active learning for deep networks: a margin  
352 based approach. *arXiv preprint arXiv:1802.09841*, 2018.
- 353 [14] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set  
354 approach. *arXiv preprint arXiv:1708.00489*, 2017.
- 355 [15] Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agar-  
356 wal. Deep batch active learning by diverse, uncertain gradient lower bounds. *arXiv preprint  
357 arXiv:1906.03671*, 2019.
- 358 [16] Kun-Peng Ning, Xun Zhao, Yu Li, and Sheng-Jun Huang. Active learning for open-set  
359 annotation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
360 Recognition*, pages 41–49, 2022.
- 361 [17] Sanjoy Dasgupta. Coarse sample complexity bounds for active learning. *Advances in neural  
362 information processing systems*, 18, 2005.
- 363 [18] David D Lewis. A sequential algorithm for training text classifiers: Corrigendum and additional  
364 data. In *Acm Sigir Forum*, volume 29, pages 13–19. ACM New York, NY, USA, 1995.
- 365 [19] Claude E Shannon. A mathematical theory of communication. *The Bell system technical  
366 journal*, 27(3):379–423, 1948.

- 367 [20] Tobias Scheffer, Christian Decomain, and Stefan Wrobel. Active hidden markov models for  
368 information extraction. In *Advances in Intelligent Data Analysis: 4th International Conference,  
369 IDA 2001 Cascais, Portugal, September 13–15, 2001 Proceedings 4*, pages 309–318. Springer,  
370 2001.
- 371 [21] Sanjoy Dasgupta. Two faces of active learning. *Theoretical computer science*, 412(19):1767–  
372 1781, 2011.
- 373 [22] Toan Tran, Thanh-Toan Do, Ian Reid, and Gustavo Carneiro. Bayesian generative active deep  
374 learning. In *International Conference on Machine Learning*, pages 6295–6304. PMLR, 2019.
- 375 [23] Keze Wang, Dongyu Zhang, Ya Li, Ruimao Zhang, and Liang Lin. Cost-effective active learning  
376 for deep image classification. *IEEE Transactions on Circuits and Systems for Video Technology*,  
377 27(12):2591–2600, 2016.
- 378 [24] Clemens-Alexander Brust, Christoph Käding, and Joachim Denzler. Active learning for deep  
379 object detection. *arXiv preprint arXiv:1809.09875*, 2018.
- 380 [25] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*,  
381 volume 4. Springer, 2006.
- 382 [26] Abhijit Bendale and Terrance E Boult. Towards open set deep networks. In *Proceedings of the  
383 IEEE conference on computer vision and pattern recognition*, pages 1563–1572, 2016.
- 384 [27] Wei Gao, Bin-Bin Yang, and Zhi-Hua Zhou. On the resistance of nearest neighbor to random  
385 noisy labels. *arXiv preprint arXiv:1607.07526*, 2016.
- 386 [28] Zhaowei Zhu, Yiwen Song, and Yang Liu. Clusterability as an alternative to anchor points  
387 when learning with noisy labels. In *International Conference on Machine Learning*, pages  
388 12912–12923. PMLR, 2021.
- 389 [29] Zhaowei Zhu, Zihao Dong, and Yang Liu. Detecting corrupted labels without training a model  
390 to predict. In *International Conference on Machine Learning*, pages 27412–27427. PMLR,  
391 2022.
- 392 [30] Kamalika Chaudhuri and Sanjoy Dasgupta. Rates of convergence for nearest neighbor classifi-  
393 cation. *Advances in Neural Information Processing Systems*, 27, 2014.
- 394 [31] Nick Rittler and Kamalika Chaudhuri. A two-stage active learning algorithm for  $k$ -nearest  
395 neighbors. *arXiv preprint arXiv:2211.10773*, 2022.
- 396 [32] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
397 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
398 models from natural language supervision. In *International conference on machine learning*,  
399 pages 8748–8763. PMLR, 2021.
- 400 [33] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
401 2009.
- 402 [34] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- 403 [35] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
404 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
405 pages 770–778, 2016.
- 406 [36] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine  
407 learning research*, 9(11), 2008.