
Active Learning for Multiple Target Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We describe and explore a novel setting of active learning (AL), where there are
2 multiple target models to be learned simultaneously. In many real applications, the
3 machine learning system is required to be deployed on diverse devices with varying
4 computational resources (e.g., workstation, mobile phone, edge devices), which
5 leads to the demand of training multiple target models on the same labeled dataset.
6 However, it is generally believed that AL is model-dependent and untransferable,
7 i.e., the data queried by one model may be less effective for training another model.
8 This phenomenon naturally raises a question “*Does there exist an AL method that*
9 *is effective for multiple target models?*”. In this paper, we answer this question by
10 theoretically analyzing the label complexity of active and passive learning under
11 the setting with multiple target models, and conclude that AL does have potential to
12 achieve better label complexity under this novel setting. Based on this insight, we
13 further propose an agnostic AL sampling strategy to **select the examples located in**
14 **the joint disagreement regions of different target models.** The experimental results
15 on the OCR benchmarks show that the proposed method can significantly surpass
16 the traditional active and passive learning methods under this challenging setting.

17 1 Introduction

18 Data labeling is expensive due to the involving of human annotator. Active learning (AL) is one of
19 the main approaches to reduce the labeling cost [27]. It evaluates the utility of the unlabeled data
20 based on the target model, and actively queries the labels from the oracle for the examples that is the
21 most beneficial to the performance improvement of the target model.

22 Existing active learning methods assume that there is only one specific target model, and try to fit
23 it with least queries. However, in many real applications, the machine learning system is required
24 to be deployed on multiple types of devices with different resource constraints [6]. For example, a
25 speech recognition software needs to support diverse machines with varying hardware efficiency,
26 ranging from high-performance workstation to the mobile-phone. Due to the different computational
27 resources, the applicable model architectures vary a lot, e.g., a deep model which is well-performed
28 on the cloud server may not be deployed on the edge device. It thus raises the demand of training
29 multiple models with different complexity to accommodate these devices.

30 Given multiple target models, how to effectively improve them with least labeled data becomes a
31 practical and challenging problem. It is generally believed that AL is usually model-dependent and
32 untransferable [21, 23, 36], i.e., the best query strategy for different target models varies a lot [37]. In
33 other words, the data queried by one model may be less effective for training another model [21].
34 These observations imply that the existing active query strategies can hardly benefit all target models
35 simultaneously, and the design of AL algorithm for multi-models can be rather difficult. A natural
36 question might be asked: “*Does there exist an active learning method which queries a set of labeled*
37 *data, such that all the target models can be effectively trained with them?*”

In this paper, we formally define the active learning for multiple target models problem, and reveal the potential improvement of AL under this novel setting. Based on this insight, we further propose an agnostic disagreement-based selection criterion. Specifically, we first define and analyze the label complexity for both active and passive learning under the setting with multiple target models. This label complexity characterizes the number of labeled examples sufficient to train an ε -good classifier with probability at least $1 - \delta$ for each target model. Moreover, we find that the label complexity of single model has a close relation to that of multiple models under the realizable case, e.g., the former provides an upper bound of the label complexity for multiple models, which also implies the potential improvement of AL under this setting. To further explore the agnostic case, we propose an active selection method DIAM (i.e., Disagreement-based AL for Multi-models) to effectively select the best examples that are beneficial to all target models. It prefers the data located in the joint disagreement regions of different models, which is expected to have higher potential to reduce the soft version space (i.e., the set of hypotheses with lower errors). Experiments are conducted on the OCR benchmarks to validate the necessity of designing active query method under this practical setting and the effectiveness of the proposed approach. The results show that the DIAM method can significantly reduce the number of queries to achieve a higher mean accuracy for multiple models compared to the traditional active and passive learning methods.

The rest of the paper is organized as follows. related work is first reviewed in the following section, then we formally define the AL for multiple target models problem and provide a general result to bridge the label complexity between single and multiple models. Then, we reveal the potential improvement of AL under this novel setting. After that, an agnostic active selection criterion is proposed and analyzed, followed by the empirical studies. And at last we conclude this work.

2 Related work

Active learning has received much attention in recent years due to the greatly increasing demands of labeling data to effectively train more complex models (e.g., deep models) [24]. One of the cores of AL is how to evaluate the potential contribution to the performance improvement of the target model for each candidate query. Most of existing criteria for active learning can be categorized into informativeness and representativeness. The informativeness-based methods [12, 34, 16] prefer the data which is near the decision boundary, and the representativeness-based methods [26, 29, 20] impose the constraints to regularize the queried data to be dissimilar with each other or conform to the latent data distribution. Many works also try to combine both criteria to obtain better performances [10, 35, 30]. Beyond these hand-crafted selection criteria, several meta-active-learning methods [17, 22, 33] are proposed to learn a generalizable query strategy across tasks. Most of the existing active learning query strategies target on improving one specific target model. They are less applicable to the multiple target models setting.

From the theoretical view, active learning theory has also been widely studied under certain conditions (e.g., binary classification, finite VC dimension) [15]. One of the interested properties of an active learning algorithm is the label complexity, which characterizes the number of queries needed to obtain an ε -good classifier with probability at least $1 - \delta$ [13]. To bound this value, disagreement coefficient [5, 4] and Shattering [14, 7] are two commonly used techniques. While most works deal with the single model setting, Balcan *et al.* [3] study the label complexity of the hypothesis space and its subclasses, which sheds light on this work. However, they mainly focus on how to construct subclasses to achieve a certain label complexity, while we aim to find an effective active learning algorithm on the given target models.

Recently, some AL studies tackle with a related problem that the target model prior can not be obtained. In this setting, they have to not only search the appropriate target model for the current task, but also avoid noneffective querying. To this end, ALMS [1] either randomly labels data to calculate the unbiased validation error for model selection, or queries by the expected error reduction to improve the models. Active-iNAS [11] considers the deep learning setting, the authors on one hand perform Neural Architecture Search (NAS) to find the appropriate model architecture, on the other hand query the examples by the searched network. Recently, Tang and Huang [31] propose a unified framework DUAL to solve this problem. They query the data that is beneficial to not only the winner model, but also the model search to identify the high potential model with least queries. All these methods try to search effective model configurations, but not improve multiple target models, which are different from our work.

3 Label Complexity of Single Model and Multiple Models

3.1 Notations and Definitions

Suppose the data is sampled from an unknown distribution \mathcal{D}_{XY} over the feature space \mathcal{X} and label space \mathcal{Y} . Given a dataset with n instances, which includes a small labeled set $\mathcal{L} = \{(\mathbf{x}_i, y_i)\}_{i=1}^{n_l}$ with n_l instances, and a large unlabeled set $\mathcal{U} = \{\mathbf{x}_i\}_{i=n_l+1}^{n_l+n_u}$ with n_u instances, where $n_l \ll n_u$ and $n = n_l + n_u$. At each iteration, the active learning method will select a batch of b examples \mathcal{Q} from \mathcal{U} for querying.

In the multiple target models setting, there are k hypothesis spaces $\{\mathcal{C}_i | i = 1, 2, \dots, k\}$ with different complexity, our goal is to actively query a set of examples to output a well-performed hypothesis \hat{h}_i from each \mathcal{C}_i . We define the true error of an hypothesis as $\text{er}(h) = \mathbb{P}_{\mathbf{x} \sim \mathcal{D}_X}(h(\mathbf{x}) \neq h^*(\mathbf{x}))$, where h^* is the target concept. The empirical error of h on \mathcal{L} is defined by $\text{er}_{\mathcal{L}}(h)$. Let $\text{Log}(a) = \max\{\ln(a), 1\}$, $\forall a > 0$, and $\mathbb{I}[\cdot]$ be the indicator function.

Here we introduce the definition of the pseudo-metric between hypotheses, which is frequently used in the subsequent proof.

Definition 1. Pseudo-metric between Hypotheses: Given marginal data distribution \mathcal{D}_X , the probability of disagreement between two classifiers h_1 and h_2 is defined as $d(h_1, h_2) = \mathbb{P}_{\mathbf{x} \sim \mathcal{D}_X}(h_1(\mathbf{x}) \neq h_2(\mathbf{x}))$.

Then we formally define the label complexity of active learning for multiple target models. Firstly, we introduce the label complexity from the common AL literature as a reference [15].

Definition 2. Label Complexity of AL: For any active learning algorithm \mathcal{A} , we say \mathcal{A} achieves label complexity Λ on the hypothesis space \mathcal{C} if, for every $\varepsilon, \delta \in (0, 1)$, every distribution \mathcal{D}_{XY} over $\mathcal{X} \times \mathcal{Y}$, and every integer $t \geq \Lambda(\varepsilon, \delta, \mathcal{D}_{XY})$, if $h_{t,\delta}$ is the classifier produced by running \mathcal{A} with budget t , then with probability at least $1 - \delta$, $\text{er}(h_{t,\delta}) - \min_{h \in \mathcal{C}} \text{er}(h) \leq \varepsilon$.

By replacing the hypothesis space in the above definition with multiple models, the label complexity for the AL with multiple target models is defined as

Definition 3. Label Complexity of AL for Multiple Target Models: Given a set of target models $T = \{\mathcal{C}_i | i = 1, 2, \dots, k\}$. For any active learning algorithm \mathcal{A} , we say \mathcal{A} achieves label complexity Λ_m for multiple target models if, for every $\varepsilon, \delta \in (0, 1)$, every distribution \mathcal{D}_{XY} over $\mathcal{X} \times \mathcal{Y}$, and every integer $t \geq \Lambda_m(\varepsilon, \delta, \mathcal{D}_{XY}, T)$, if $\{h_i^{t,\delta} \in \mathcal{C}_i | i = 1, \dots, k\}$ is the classifiers produced by running \mathcal{A} with budget t , then with probability at least $1 - \delta$, $\text{er}(h_i^{t,\delta}) - \min_{h_i \in \mathcal{C}_i} \text{er}(h_i) \leq \varepsilon, \forall i = 1, \dots, k$.

3.2 A Bridge of Label Complexity between Single and Multiple Models in Realizable Case

Trivially, the label complexity for multiple models must has $\Lambda_m \leq \sum_i \Lambda_i$. (applying the AL algorithm \mathcal{A} on each of the target model i to get the result.) However, such result mainly claims that the AL tends to be non-effective under this setting (cf. Theorem 2 for the label complexity of multiple models of passive learning). To break this curse, the following theorem is provided to show that, we can expect a much better performance for AL under this setting. It generally says that if a learning method has label complexity Λ on hypothesis space \mathcal{C} , it also has the ability to output good classifiers for any subsets of the hypothesis space, i.e., after querying at most t examples to output ε -good classifiers with probability at least $1 - \delta$ for any sequence of subsets $\mathcal{C}_1, \mathcal{C}_2, \dots$ with $\bigcup_{i=1}^{\infty} \mathcal{C}_i = \mathcal{C}$.

Theorem 1. Considering binary classification tasks and realizable case, given hypothesis space \mathcal{C} , assume that active learning algorithm \mathcal{A} achieves label complexity Λ on \mathcal{C} . Then, for any sequence of subsets $\mathcal{C}_1, \mathcal{C}_2, \dots$ with $\bigcup_{i=1}^{\infty} \mathcal{C}_i = \mathcal{C}$, there exists a global active learning algorithm \mathcal{A}' which achieves the label complexity Λ_m for multiple target models $T = \{\mathcal{C}_i | i = 1, 2, \dots\}$ such that $\Lambda_m(\varepsilon, \delta, \mathcal{D}_{XY}, T) = \Lambda(\varepsilon/2, \delta, \mathcal{D}_{XY})$.

Proof. Define an algorithm \mathcal{A}' that can output the required classifier $\hat{h}_i \in \mathcal{C}_i$ for any \mathcal{C}_i as follows. First, run the algorithm \mathcal{A} on $(\mathcal{C}, \mathcal{D}_{XY})$ to query $t \geq \Lambda(\varepsilon/2, \delta, \mathcal{D}_{XY})$ labels and output a classifier h_A . Then, for any given \mathcal{C}_i , output the classifier $\hat{h}_i \in \mathcal{C}_i$ such that $\hat{h}_i = \min_{h_i \in \mathcal{C}_i} d(h_i, h_A)$. Next, we prove that $\text{er}(\hat{h}_i) - \min_{h_i \in \mathcal{C}_i} \text{er}(h_i) \leq \varepsilon$ holds with probability at least $1 - \delta$.

142 To bound $\text{er}(\hat{h}_i)$, it is equivalent to bound $d(\hat{h}_i, h^*)$ by Definition 1. Let $h_i^* = \arg \min_{h_i \in \mathcal{C}_i} \text{er}(h_i)$.
 143 It is easy to verify that, $d(\cdot)$ satisfies triangle inequality in binary classification problems, i.e.,

$$d(\hat{h}_i, h^*) \leq d(\hat{h}_i, h_A) + d(h_A, h^*). \quad (1)$$

144 For the $d(\hat{h}_i, h_A)$, we know that $\hat{h}_i = \min_{h_i \in \mathcal{C}_i} d(h_i, h_A)$, which means

$$d(\hat{h}_i, h_A) \leq d(h_i^*, h_A). \quad (2)$$

145 Again, by the triangle inequality, we have

$$d(h_i^*, h_A) \leq d(h_i^*, h^*) + d(h^*, h_A). \quad (3)$$

146 Combining Eq. (1)(2)(3), we have

$$d(\hat{h}_i, h^*) \leq d(h_i^*, h^*) + 2d(h_A, h^*). \quad (4)$$

147 Since $d(h_A, h^*)$ is bounded by $\varepsilon/2$ with probability at least $1 - \delta$ according to the definition, we can
 148 get $\text{er}(\hat{h}_i) - \min_{h_i \in \mathcal{C}_i} \text{er}(h_i) \leq \varepsilon$ holds with probability at least $1 - \delta$. \square

149 Theorem 1 says that if we can run an active learning method to obtain a classifier h_A such that
 150 $\text{er}(h_A) \leq \varepsilon/2$ in \mathcal{C} , then we can obtain ε -good classifier \hat{h}_i with probability at least $1 - \delta$ for any
 151 subset \mathcal{C}_i by $\hat{h}_i = \min_{h_i \in \mathcal{C}_i} d(h_i, h_A)$. This result provides a general guarantee that if an algorithm
 152 can achieve a label complexity on the combined hypothesis space of different models, it also can
 153 achieve a bounded label complexity on these models (i.e., the label complexity for multiple models).
 154 It can be served as a baseline of AL with multi-models setting.

155 4 Potential Improvements of Active over Passive

156 Although Theorem 1 bridges the traditional label complexity to that of multiple models setting, it
 157 does not reveal the improvement of active over passive learning. Next, we will discuss this topic.

158 4.1 Label Complexity of Passive Learning for Multiple Target Models

159 According to the theoretical analysis of the empirical risk minimization (ERM) [15], we know that

160 **Lemma 1.** *Considering the binary classification, given the hypothesis space \mathcal{C} with VC dimension d .
 161 The passive learning algorithm $\text{ERM}(\mathcal{C}, \cdot)$ achieves a label complexity Λ such that, for any \mathcal{D}_{XY} in
 162 the realizable case, $\forall \varepsilon, \delta \in (0, 1)$,*

$$\Lambda(\varepsilon, \delta, \mathcal{D}_{XY}) \lesssim \left(\frac{1}{\varepsilon} \right) (d \text{Log}(\theta(\varepsilon)) + \text{Log}(1/\delta)). \quad (5)$$

163 For the agnostic case, define $\nu = \min_{h \in \mathcal{C}} \text{er}(h)$, $\text{ERM}(\mathcal{C}, \cdot)$ achieves a label complexity Λ such that,
 164

$$\Lambda(\nu + \varepsilon, \delta, \mathcal{D}_{XY}) \lesssim \left(\frac{\nu + \varepsilon}{\varepsilon^2} \right) (d \text{Log}(\theta(\nu + \varepsilon)) + \text{Log}(1/\delta)), \quad (6)$$

165 where $\theta(\cdot)$ is the disagreement coefficient which is formally defined as

Definition 4. Disagreement Coefficient: For any $r_0 \geq 0$ and classifier h , define the disagreement coefficient of h with respect to \mathcal{C} on \mathcal{D}_{XY} as

$$\theta_h^{\mathcal{C}}(r_0) = \sup_{r > r_0} \frac{\mathbb{P}(\text{DIS}(\text{B}_{\mathcal{C}}(h, r)))}{r} \vee 1.$$

166 Where \vee is the max operator, $\text{DIS}(\mathcal{H}) = \{\mathbf{x} \in \mathcal{X} \mid \exists h, h' \in \mathcal{H}, \text{ s.t. } h(\mathbf{x}) \neq h'(\mathbf{x})\}$, $\text{B}_{\mathcal{H}}(h, r) =$
 167 $\{g \in \mathcal{H} \mid d(h, g) \leq r\}$ for a given set of hypotheses \mathcal{H} . Specifically, if $\mathcal{H} = \mathcal{C}$, abbreviate $\text{B}_{\mathcal{C}}(h, r) =$
 168 $\text{B}(h, r)$ and $\theta_h^{\mathcal{C}} = \theta_h$. Suppose $f^* = \arg \min_{h \in \mathcal{C}} \text{er}(h)$, We further abbreviate $\theta_{f^*}(r_0) = \theta(r_0)$.

169 This value roughly characterizes the behavior of the size of disagreement region $\text{DIS}(\cdot)$ as a function
 170 of the hypotheses within a radius r around the classifier h . With Lemma 1, we can easily obtain the
 171 following result

Theorem 2. Denote by Λ_i the label complexity that the passive learning achieved on the i -th target model \mathcal{C}_i . The passive learning will achieves a label complexity Λ_m for multiple target models such that $\Lambda_m = \max_i \Lambda_i$ on any given target models $T = \{\mathcal{C}_i | i = 1, 2, \dots, k\}$.

Proof. Since the data is randomly sampled, if $t \geq \max_i \Lambda_i(\varepsilon, \delta, \mathcal{D}_{XY})$ examples are queried, according to the definition of label complexity, the passive learning will output an ε -good classifier with probability at least $1 - \delta$ for all target models. \square

This result shows that the random sampling can be a very competitive method for the multiple target models setting. However, we note that the target concept h^* will usually not be included by every hypothesis space \mathcal{C}_i , thus its label complexity will heavily depends on the property of the given models, i.e., the value of $\min_{h \in \mathcal{C}_i} \text{er}(h)$. In this way, we can expect an improvement of Theorem 1 (though not always) by applying an existing active learning algorithm on the combined hypothesis space \mathcal{C} in the realizable case.

4.2 Potential Improvement of Active Learning for Multiple Target Models

To show the potential of AL under this setting, we take the CAL method [9] as an example, which is a representative and well-analyzed approach in the active learning literature [15]. CAL queries the examples from the disagreement region of a set of consistent hypotheses, i.e., $\text{DIS}(V) = \{\mathbf{x} \in \mathcal{X} \mid \exists h, h' \in V \text{ s.t. } h(\mathbf{x}) \neq h'(\mathbf{x})\}$, where $V = \{h \in \mathcal{C} \mid h(\mathbf{x}) = h^*(\mathbf{x}), \forall \mathbf{x} \in \mathcal{L}\}$. It achieves the label complexity $O(\theta(\varepsilon) \log(1/\varepsilon) \log(\theta(\varepsilon) \log(1/\varepsilon)))$ for the realizable case. According to Theorem 1, it will have the following label complexity for the multiple target models

Corollary 1. Given target models $T = \{\mathcal{C}_i | i = 1, 2, \dots, k\}$, define $\mathcal{C} = \bigcup_{i=1}^k \mathcal{C}_i$. Suppose \mathcal{C} has VC dimension $d < \infty$. CAL achieves a label complexity Λ_m for multiple target models such that, for \mathcal{D}_{XY} in the realizable case, for any $\forall \varepsilon, \delta \in (0, 1)$,

$$\Lambda_m(\varepsilon, \delta, \mathcal{D}_{XY}, T) \leq \theta_{h^*}^{\mathcal{C}}(\varepsilon/2) \text{Log}(2/\varepsilon) \left(d \text{Log}(\theta_{h^*}^{\mathcal{C}}(\varepsilon/2)) + \text{Log}\left(\frac{\text{Log}(2/\varepsilon)}{\delta}\right) \right). \quad (7)$$

Proof. By combining the label complexity of CAL for single model from [15] and Theorem 1, we can get the result. \square

To reveal the potential improvement, assume that there exists a target model \mathcal{C}_m such that $\min_{h \in \mathcal{C}_m} \text{er}(h) > \varepsilon$. Then according to Lemma 1 and Theorem 2, the label complexity of passive learning for multiple target models $\Lambda_m(\varepsilon, \delta, \mathcal{D}_{XY}, T)$ is $\Omega(2/\varepsilon)$. On the other side, Corollary 1 shows that, CAL has a label complexity $\Omega(\text{Log}(2/\varepsilon))$, which implies the potential of the improvement of active learning under this setting. We leave the guarantee of strict improvement of AL under this setting an interesting future work. Note that CAL only works in the realizable case, next we further study the agnostic case.

5 An Agnostic Disagreement-based AL method for Multiple Models

Define the set V_i for each \mathcal{C}_i as $\{h \in \mathcal{C}_i \mid h(\mathbf{x}) = h^*(\mathbf{x}), \forall \mathbf{x} \in \mathcal{L}\}$, we propose to query the examples located in the joint disagreement regions for all $\mathcal{C}_i, \forall i = 1, 2, \dots$, i.e., $\text{DIS}(V_1) \cap \text{DIS}(V_2) \cap \dots$. Intuitively, we know that V_i must be a subset of V , if such data exists, we can expect it has higher potential to reduce V . This statement can be simply implied by the Bayesian formula.

Proposition 1. Considering binary classification problem, Let $V_+(\mathbf{x}) = \{h \in V \mid h(\mathbf{x}) = +1\}$, $V_-(\mathbf{x}) = \{h \in V \mid h(\mathbf{x}) = -1\}$, and $\lambda(\mathbf{x}) = \frac{|V_+(\mathbf{x})|}{|V_-(\mathbf{x})|}$, where $|\cdot|$ is the number of elements of a set. The ideal case is to query the \mathbf{x} which has $\lambda(\mathbf{x}) = 1$. Given any sequences of subset V_1, V_2, \dots, V_k randomly sampled from V , define the event $E_{\mathbf{x}}$ that data \mathbf{x} falls into $\text{DIS}(V_1) \cap \text{DIS}(V_2) \cap \dots \cap \text{DIS}(V_k)$. According to the Bayesian formula, we have

$$\begin{aligned} \mathbb{P}(\lambda(\mathbf{x}) = 1 | E_{\mathbf{x}}) &= \frac{\mathbb{P}(E_{\mathbf{x}} | \lambda(\mathbf{x}) = 1) \mathbb{P}(\lambda(\mathbf{x}) = 1)}{\mathbb{P}(E_{\mathbf{x}})} \\ &\geq \mathbb{P}(\lambda(\mathbf{x}) = 1). \end{aligned} \quad (8)$$

Algorithm 1 The DIAM-online Algorithm

Initialize: hyperparameter q , constants σ_i ;
 $m \leftarrow 0, \hat{V}_i \leftarrow \mathcal{C}_i, \forall i = 1, \dots, k$.
Output: Any $h \in \hat{V}_i, \forall i = 1, \dots, k$.

```

1: while Labeling budget is not run out do
2:    $m \leftarrow m + 1$ 
3:   Request an unlabeled data  $\mathbf{x}_m$ 
4:   if  $\sum_i \mathbb{I}[\mathbf{x}_m \in \text{DIS}(\hat{V}_i)] \geq q$  then
5:     Query:  $\mathcal{L} \leftarrow \mathcal{L} \cup \{(\mathbf{x}_m, h^*(\mathbf{x}_m))\}$ 
6:   end if
7:   if  $\log_2 m \in \mathbb{N}$  then
8:      $\hat{V}_i \leftarrow \{h \in \hat{V}_i \mid \text{er}_{\mathcal{L}}(h) - \min_{g \in \hat{V}_i} \text{er}_{\mathcal{L}}(g) \leq \sigma_i\}, \forall i = 1, \dots, k$ .
9:   end if
10: end while

```

Algorithm 2 The DIAM-pool Algorithm

Initialize: labeled set \mathcal{L} , unlabeled set \mathcal{U} , hyper-parameters $\hat{\sigma}_i, \hat{V}_i \leftarrow \mathcal{C}_i, \forall i = 1, \dots, k$.
Output: $\{\hat{h}_i \mid i = 1, \dots, k\}$.

```

1: while Labeling budget is not run out do
2:    $\mathbf{x}^* = \arg \max_{\mathbf{x} \in \mathcal{U}} \sum_i \mathbb{I}[\mathbf{x} \in \text{DIS}(\hat{V}_i)]$ 
3:   Query  $\mathbf{x}^*$  from the oracle:  $\mathcal{L} \leftarrow \mathcal{L} \cup \{(\mathbf{x}^*, h^*(\mathbf{x}^*))\}$ 
4:    $\mathcal{U} \leftarrow \mathcal{U} \setminus \{\mathbf{x}^*\}$ 
5:   for  $i = 1, \dots, k$  do
6:      $\hat{h}_i \leftarrow \min_{g \in \hat{V}_i} \text{er}_{\mathcal{L}}(g)$ 
7:      $\hat{V}_i \leftarrow \{h \in \hat{V}_i \mid (\text{er}_{\mathcal{L}}(h) - \hat{h}_i) \leq \hat{\sigma}_i\}$ 
8:   end for
9: end while

```

213 *Proof.* Since each V_i is randomly sampled from V , $\mathbb{P}(E_{\mathbf{x}})$ will reach its maximum value when
 214 $\lambda(\mathbf{x}) = 1$, thus we have $\mathbb{P}(E_{\mathbf{x}} \mid \lambda(\mathbf{x}) = 1) / \mathbb{P}(E_{\mathbf{x}}) \geq 1$, which leads to the conclusion. \square

215 Following this principle, we would like to query the examples located in the joint disagreement
 216 regions of $\mathcal{C}_i, \forall i = 1, 2, \dots, k$. However, since we have multiple target models, the target concept
 217 h^* might not be included by every \mathcal{C}_i in practice, which turns the learning problem to the agnostic
 218 setting. Inspired by the RobustCAL method [2], which is a disagreement-based AL algorithm for the
 219 agnostic setting, we propose DIAM (i.e., DISagreement-based AL for Multi-models) query strategy
 220 for the multiple target models problem. Note that we define a new form of V as \hat{V} to tackle with the
 221 noisy setting, i.e., $\hat{V} = \{h \in \mathcal{C} \mid \text{er}_{\mathcal{L}}(h) - \min_{g \in \mathcal{C}} \text{er}_{\mathcal{L}}(g) \leq \sigma\}$, where σ is a constant. To simplify
 222 the theoretical analysis, we first propose an online version of DIAM, then we define the DIAM
 223 method for the pool-based setting and empirically validate its effectiveness. They are summarized at
 224 Algorithm 1 and 2, respectively.

225 Now let us analyze the DIAM method. Since we are considering the agnostic setting, it is necessary
 226 to model the noise. Here we employ the commonly used Tsybakov noise condition [32].

Condition 1. [32, Tsybakov noise] For some $a \in [1, \infty)$ and $\alpha \in [0, 1]$, assume that f^* achieves
 $\inf_{h \in \mathcal{C}} \text{er}(h)$, for every $h \in \mathcal{C}$,

$$\mathbb{P}(\mathbf{x} : h(\mathbf{x}) \neq f^*(\mathbf{x})) \leq a (\text{er}(h) - \text{er}(f^*))^\alpha.$$

227 We assume that there exists pairs of a_i and α_i for each target model \mathcal{C}_i . By further taking the constants
 228 σ_i in DIAM-online algorithm as the same form in the RobustCAL method [15]. Considering a
 229 conservative situation that the hyperparameter $q = 1$, we can have the following result. The proof is
 230 deferred to the appendix.

231 **Theorem 3.** Considering agnostic setting and binary classification tasks. Given a set of target
 232 models $T = \{\mathcal{C}_i \mid i = 1, 2, \dots, k\}$ who have VC dimensions $d_i < \infty$ and meet Condition 1. Let
 233 $h_i^* = \arg \min_{h_i \in \mathcal{C}_i} \text{er}(h_i)$ and $\nu_i = \text{er}(h_i^*)$. For any $\varepsilon, \delta \in (0, 1)$, if $q = 1$, DIAM-online algorithm
 234 achieves a label complexity $\Lambda_m^{\varepsilon, \delta} = \Lambda_m(\varepsilon, \delta, \mathcal{D}_{XY}, T)$ for multiple target models such that, for a_i
 235 and α_i as in Condition 1, for any \mathcal{D}_{XY} ,

$$\Lambda_m^{\varepsilon, \delta} \leq \sum_{i=1}^k a_i^2 \theta_{h_i^*}^{\mathcal{C}_i}(a_i \varepsilon^{\alpha_i}) \varepsilon^{2\alpha_i - 2} \left(d_i \text{Log} \left(\theta_{h_i^*}^{\mathcal{C}_i}(a_i \varepsilon^{\alpha_i}) \right) + \text{Log} \left(\frac{\text{Log}(a_i / \varepsilon)}{\delta} \right) \right) \text{Log} \left(\frac{1}{\varepsilon} \right), \quad (9)$$

236 and furthermore,

$$\Lambda_m^{\varepsilon, \delta} \leq \sum_{i=1}^k \theta_{h_i^*}^{\mathcal{C}_i}(\nu_i + \varepsilon) \left(\frac{\nu_i^2}{\varepsilon^2} + \text{Log} \left(\frac{1}{\varepsilon} \right) \right) \left(d_i \text{Log}(\theta_{h_i^*}^{\mathcal{C}_i}(\nu_i + \varepsilon)) + \text{Log} \left(\frac{\text{Log}(1/\varepsilon)}{\delta} \right) \right). \quad (10)$$

Theorem 3 provides an upper bound of the label complexity of the DIAM-online method when $q = 1$. It considers a general situation with arbitrary target models and data distributions, even the unlabeled data will never fall into the joint disagreement regions. However, one may be more interested in the situation that if we can always query the \mathbf{x} such that $\sum_i \mathbb{I}[\mathbf{x} \in \text{DIS}(\hat{V}_i)] = k$. Next, we prove that if such ideal situation exists, DIAM-online achieves a better label complexity than applying CAL on the multiple target models setting under the realizable case.

Theorem 4. *Considering binary classification tasks. Given a set of target models $T = \{C_i | i = 1, 2, \dots, k\}$ who have VC dimensions $d_i < \infty$ and meet Condition 1, such that, if a data point falls into the disagreement region of any C_i , it also falls into the disagreement regions of the others $\{C_j | j \neq i, j = 1, 2, \dots, k\}$. Assume the m -th target model achieves the highest label complexity. Let $h_m^* = \arg \min_{h_m \in C_m} \text{er}(h_m)$ and $\nu_m = \text{er}(h_m^*)$. For any $\delta \in (0, 1)$, $\varepsilon \in (0, 1/e)$, $h^* \in C$ where $C = \bigcup_{i=1}^{\infty} C_i$ with VC dimension $d < \infty$. If $\nu_m \leq \frac{\ln 2}{2} \varepsilon$, DIAM-online achieves a better label complexity for multiple target models than that of applying CAL method on C . More specifically, the following inequality holds:*

$$\begin{aligned} & \theta_{h^*}^C(\varepsilon/2) \ln(2/\varepsilon) \left(d \ln(\theta_{h^*}^C(\varepsilon/2)) + \ln \left(\frac{\ln(2/\varepsilon)}{\delta} \right) \right) > \\ & \theta_{h_m^*}^{C_m}(\nu_m + \varepsilon) \left(\frac{\nu_m^2}{\varepsilon^2} + \ln \left(\frac{1}{\varepsilon} \right) \right) \left(d_m \ln(\theta_{h_m^*}^{C_m}(\nu_m + \varepsilon)) + \ln \left(\frac{\ln(1/\varepsilon)}{\delta} \right) \right). \end{aligned} \quad (11)$$

The key of the proof is comparing the disagreement coefficients defined on different functions and hypothesis spaces, i.e., $\theta_{h_m^*}^{C_m}$ and $\theta_{h^*}^C$. We defer the proof to the appendix. Although Theorem 4 holds with somewhat strict conditions, we note that Theorem 1 only works in the realizable case, i.e., $h^* \in C$, while DIAM does not require this constraint, it is proposed for the agnostic setting. Next, we discuss how to implement DIAM in the real applications for deep models.

It is generally believed that finding the disagreed pair of classifiers from a set of hypotheses for a given \mathbf{x} is non-trivial. Most existing methods randomly sample functions from the hypothesis space for validation, or turn to select the data close to the decision boundary (e.g., uncertainty), which can be expensive or inaccurate. This problem becomes more prohibitive to the deep models.

To efficiently estimate the disagreement regions for the neural networks, we propose to exploit the predictions of the unlabeled data during the later epochs in the training phase, typically after the network converging. Recall the definition of disagreement region $\text{DIS}(\hat{V}_i)$, we should firstly find the hypotheses that are basically consistent with the labeled data, then validate whether there exists a pair of hypotheses that disagree on the given unlabeled data. For the first characteristic, the models on the later epochs, i.e., has smaller training errors, can represent the well-learned hypotheses. For the second characteristic, if there exists models i, j from the later epochs such that the model trained at epoch i has inconsistent prediction with the model trained at epoch j on the unlabeled data \mathbf{x} , we can say that the example \mathbf{x} falls into $\text{DIS}(\hat{V}_i)$. We also note that the query batch size in deep learning is usually large, to avoid overmuch information redundancy, we heuristically sort the unlabeled data according to the active selection scores, and randomly select a batch of examples from the top-rated candidates, e.g., 5 times the size of the batch size.

6 Experiment

6.1 Empirical Settings

To construct the multiple target models scenario, we introduce the results of a recent NAS method OFA [6], which tries to efficiently search model architectures for different devices by training only one super-net. They report the searched effective architectures that meet the hardware constraints of various machines on the GitHub¹, which is well suited to our problem setting. Specifically, we take 12 specialized model architectures with different prediction accuracies and speeds that target on Samsung S7 Edge, Samsung Note8 and Samsung Note10 as our target models. More details are introduced in the appendix.

We compare the following query strategies in our experiments.

¹<https://github.com/mit-han-lab/once-for-all>. It is under the MIT license.

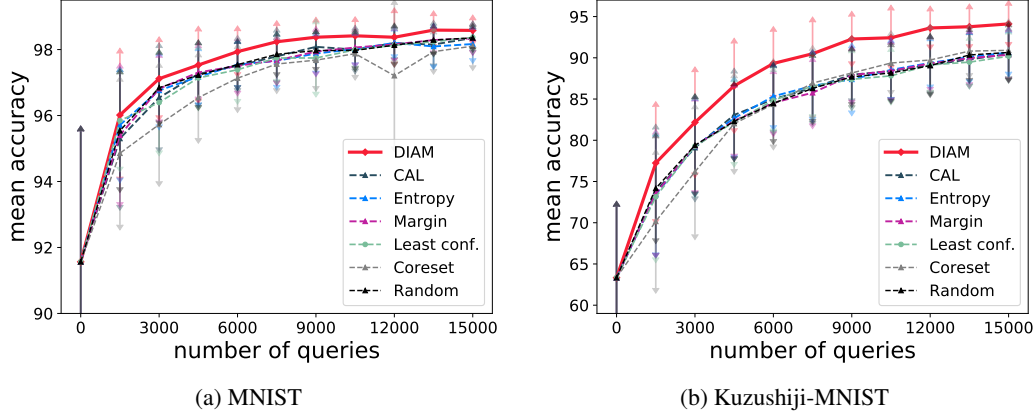


Figure 1: The learning curves with the mean accuracy of the target models of the compared methods. The error bars indicate the standard deviation of the performances of target models.

- DIAM: The proposed method of this paper, which queries the data located in the joint disagreement regions of multiple target models.
- CAL [9]: Query the data falls into the disagreement region of any target models. It has a bounded label complexity for the multiple target models setting according to Theorem 1.
- Entropy [19]: Query the data with the highest prediction entropies. We take the mean entropies calculated by all target models to support the novel problem setting.
- Least Confidence [28]: Query the data with the least prediction confidence. We take the mean values calculated by all target models to support the novel problem setting.
- Margin [25]: Query the data with the minimum prediction margin. We take the mean margin values calculated by all target models to support the novel problem setting.
- Coreset [26]: Query the most representative data. The distance is calculated by the features extracted by a pretrained MobileNetV3, which is the super-net in OFA [6].
- Random: Query data randomly. Note that this is a highly competitive baseline according to Theorem 2.

Since Optical Character Recognition (OCR) is one of the representative machine learning systems that are required to be deployed on diverse devices, two commonly used hand-writing characters classification benchmarks are employed in our experiments, i.e., the MNIST [18] and Kuzushiji-MNIST [8] datasets. They are under the CC BY-SA 3.0 and CC BY-SA 4.0 licenses, respectively. Here we consider the prevalent pool-based active learning setting. Specifically, we randomly take 3,000 training data as our initially labeled data, and the rest as the unlabeled pool. At each iteration, the compared sampling methods will select 1,500 unlabeled examples for querying, then re-train the models. The mean and standard deviation of the accuracies of multiple target models are reported. Note that more results can be found in the appendix.

For the model training, We mainly follow the training configs of OFA. However, since the initially labeled data is limited, a small number of training epochs is taken to avoid over-fitting. Specifically, we employ the pretrained weights on the image-net dataset for initialization, then finetune 20 epochs on the labeled data.

6.2 Results

We report the trend of mean accuracy of multiple target models with the number of queries increasing in Fig. 1. The error bars indicate the standard deviation of the performances. First of all, the high deviation of the performances of the initial point shows the diversity of the target models, which symbolizes the practicability and difficulty of the experimental settings. It is conceivable that different target models will have various preferences of training data due to the diverse architectures. Under this challenging setting, it can be observed from the figure that our method can significantly outperform the traditional active and passive learning methods. It shows a great potential of improvements

Table 1: The mean of the learning curves and the mean of standard deviation values with different numbers of target models on the OCR benchmarks achieved by the compared methods (mean accuracy \pm mean standard deviation). The best performance is highlighted in boldface.

Methods	Number of Target Models				
	2	4	6	8	12
MNIST					
DIAM	98.16 \pm 0.13	97.29 \pm 0.99	97.55 \pm 0.85	97.34 \pm 1.09	97.34 \pm 1.04
CAL	97.79 \pm 0.14	97.04 \pm 0.92	97.24 \pm 0.89	96.95 \pm 1.07	96.98 \pm 1.10
Entropy	97.83 \pm 0.10	96.94 \pm 1.01	97.15 \pm 0.98	96.92 \pm 1.06	96.98 \pm 1.00
Margin	97.79 \pm 0.13	96.94 \pm 1.02	97.19 \pm 0.96	96.81 \pm 1.19	97.00 \pm 1.02
Least conf.	97.84 \pm 0.11	96.89 \pm 1.02	97.23 \pm 0.92	96.88 \pm 1.05	96.96 \pm 1.07
Coreset	97.64 \pm 0.13	96.69 \pm 1.07	97.03 \pm 0.97	96.36 \pm 1.82	96.56 \pm 1.40
Random	97.81 \pm 0.12	96.93 \pm 0.97	97.21 \pm 0.94	96.83 \pm 1.12	97.03 \pm 0.99
Kuzushiji-MNIST					
DIAM	90.38 \pm 0.21	85.76 \pm 4.69	86.91 \pm 4.38	86.23 \pm 4.68	86.85 \pm 4.25
CAL	87.06 \pm 0.34	83.61 \pm 4.29	84.70 \pm 3.88	83.40 \pm 4.32	83.31 \pm 4.53
Entropy	87.09 \pm 0.34	83.22 \pm 4.16	84.39 \pm 3.85	83.28 \pm 4.28	83.33 \pm 4.33
Margin	86.91 \pm 0.35	83.20 \pm 4.10	84.31 \pm 4.03	83.11 \pm 4.37	83.16 \pm 4.29
Least conf.	86.71 \pm 0.26	83.38 \pm 4.25	84.36 \pm 3.71	83.20 \pm 4.42	83.04 \pm 4.33
Coreset	87.49 \pm 0.36	82.97 \pm 5.16	84.80 \pm 4.58	83.00 \pm 4.93	82.91 \pm 5.03
Random	87.34 \pm 0.31	82.97 \pm 4.38	84.22 \pm 3.98	83.02 \pm 4.19	83.26 \pm 4.36

over the random sampling, which is a very competitive baseline according to Theorem 2. This result sufficiently reveals the effectiveness of DIAM and the necessity of designing active query method under this practical setting. The uncertainty-based methods, i.e., entropy, least confidence and margin, achieve comparable performances with random sampling. These results meet our expectation. Because traditional AL methods are usually model-dependent, i.e., the data queried by one model may be less effective for training another model. By taking the mean uncertainty scores of diverse target models, the data selection may tend to be non-informative. The coreset method is less stable than random. We note that coreset is still a model-based selection method in deep learning. Because the features of the data will be optimized along with the training procedures. Thus it may also suffer from the model dependence problem.

6.3 Study on Different Numbers of Target Models

We further explore the influence of the number of target models to the data selection methods. The mean of the learning curves and the mean of standard deviations are reported in Table 1. The results show that our method can consistently outperform the other compared methods, which demonstrate its robustness to the number of models. This property also denotes that the DIAM method has the potential to tackle with more challenging situations, i.e., improving sufficient numbers of target models simultaneously. It is essential to the machine learning systems which have a wide range of applications. The performances of the other compared methods have similar trends with more target models setting. It again verifies that the traditional AL methods are usually model-dependent, and emphasizes the necessity of designing novel selection approaches under this practical setting.

7 Conclusion

In this paper, we propose to study active learning in a novel setting, where the task is to select and label the most useful examples that are beneficial to multiple target models. We firstly analyze the label complexity of active and passive learning to reveal the potential improvement of AL under this novel setting. Based on this insight, we further propose an active selection criterion DIAM that prefers the data located in the joint disagreement regions of different target models. Empirical studies on the OCR benchmarks, which is one of the representative applications that are required to accommodate different devices, show the effectiveness of the proposed method. In the future, we will tackle with more complex and important learning tasks (e.g., face recognition, object detection), and design effective query strategies under the multiple target models setting.

References

- [1] Alnur Ali, Rich Caruana, and Ashish Kapoor. Active learning with model selection. In *AAAI Conference on Artificial Intelligence*, pages 1673–1679, 2014.
- [2] Maria-Florina Balcan, Alina Beygelzimer, and John Langford. Agnostic active learning. *Journal of Computer and System Sciences*, 75(1):78–89, 2009.
- [3] Maria-Florina Balcan, Steve Hanneke, and Jennifer Wortman Vaughan. The true sample complexity of active learning. *Machine Learning*, 80(2):111–139, 2010.
- [4] Alina Beygelzimer, Daniel J. Hsu, John Langford, and Chicheng Zhang. Search improves label for active learning. In *Advances in Neural Information Processing Systems*, pages 3342–3350, 2016.
- [5] Alina Beygelzimer, Daniel J. Hsu, John Langford, and Tong Zhang. Agnostic active learning without constraints. In *Advances in Neural Information Processing Systems*, pages 199–207, 2010.
- [6] Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han. Once-for-all: Train one network and specialize it for efficient deployment. In *International Conference on Learning Representations*, 2020.
- [7] Xiaofeng Cao and Ivor W Tsang. Shattering distribution for active learning. *IEEE Transactions on Neural Networks and Learning Systems*, 33(1):215–228, 2022.
- [8] Tarin Clanuwat, Mikel Bober-Irizar, Asanobu Kitamoto, Alex Lamb, Kazuaki Yamamoto, and David Ha. Deep learning for classical japanese literature, 2018.
- [9] David Cohn, Les Atlas, and Richard Ladner. Improving generalization with active learning. *Machine Learning*, 15(2):201–221, 1994.
- [10] Bo Du, Zengmao Wang, Lefei Zhang, Liangpei Zhang, Wei Liu, Jialie Shen, and Dacheng Tao. Exploring representativeness and informativeness for active learning. *IEEE Transactions on Cybernetics*, 47(1):14–26, 2015.
- [11] Yonatan Geifman and Ran El-Yaniv. Deep active learning with a neural architecture search. In *Advances in Neural Information Processing Systems*, pages 5976–5986, 2019.
- [12] Ran Gilad-Bachrach, Amir Navot, and Naftali Tishby. Query by committee made real. In *Advances in Neural Information Processing Systems*, pages 443–450, 2005.
- [13] Steve Hanneke. A bound on the label complexity of agnostic active learning. In *International Conference on Machine Learning*, pages 353–360, 2007.
- [14] Steve Hanneke. Activized learning: Transforming passive to active with improved label complexity. *The Journal of Machine Learning Research*, 13(1):1469–1587, 2012.
- [15] Steve Hanneke. Theory of active learning. *Foundations and Trends in Machine Learning*, 7(2-3), 2014.
- [16] Andreas Kirsch, Joost van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning. In *Advances in Neural Information Processing Systems*, pages 7026–7037, 2019.
- [17] Ksenia Konyushkova, Raphael Sznitman, and Pascal Fua. Learning active learning from data. In *Advances in Neural Information Processing Systems*, pages 4225–4235, 2017.
- [18] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [19] David D Lewis and Jason Catlett. Heterogeneous uncertainty sampling for supervised learning. In *Machine Learning: Proceedings of the 11th International Conference*, pages 148–156. Elsevier, 1994.
- [20] Changsheng Li, Handong Ma, Zhao Kang, Ye Yuan, Xiao-Yu Zhang, and Guoren Wang. On deep unsupervised active learning. In *International Joint Conferences on Artificial Intelligence*, pages 2626–2632, 2020.
- [21] David Lowell, Zachary C. Lipton, and Byron C. Wallace. Practical obstacles to deploying active learning. In *EMNLP-IJCNLP*, pages 21–30, 2019.
- [22] Kunkun Pang, Mingzhi Dong, Yang Wu, and Timothy Hospedales. Meta-learning transferable active learning policies by deep reinforcement learning. *arXiv preprint arXiv:1806.04798*, 2018.

- [23] Davi Pereira-Santos, Ricardo Bastos Cavalcante Prudêncio, and André CPLF de Carvalho. Empirical investigation of active learning strategies. *Neurocomputing*, 326:15–27, 2019.
- [24] Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Xiaojiang Chen, and Xin Wang. A survey of deep active learning. *arXiv preprint arXiv:2009.00236*, 2020.
- [25] Tobias Scheffer, Christian Decomain, and Stefan Wrobel. Active hidden markov models for information extraction. In *International Symposium on Intelligent Data Analysis*, pages 309–318. Springer, 2001.
- [26] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. In *International Conference on Learning Representations*, 2018.
- [27] Burr Settles. Active learning literature survey. Technical report, University of Wisconsin-Madison, 2009.
- [28] Burr Settles and Mark Craven. An analysis of active learning strategies for sequence labeling tasks. In *Conference on Empirical Methods in Natural Language Processing*, pages 1070–1079, 2008.
- [29] Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell. Variational adversarial active learning. In *International Conference on Computer Vision*, pages 5972–5981, 2019.
- [30] Ying-Peng Tang and Sheng-Jun Huang. Self-paced active learning: Query the right thing at the right time. In *AAAI Conference on Artificial Intelligence*, volume 33, pages 5117–5124, 2019.
- [31] Ying-Peng Tang and Sheng-Jun Huang. Dual active learning for both model and data selection. In *International Joint Conference on Artificial Intelligence*, pages 3052–3058, 2021.
- [32] Alexander B Tsybakov. Optimal aggregation of classifiers in statistical learning. *The Annals of Statistics*, 32(1):135–166, 2004.
- [33] Thuy-Trang Vu, Ming Liu, Dinh Phung, and Gholamreza Haffari. Learning how to active learn by dreaming. In *Annual Meeting of the Association for Computational Linguistics*, pages 4091–4101, 2019.
- [34] Yifan Yan and Sheng-Jun Huang. Cost-effective active learning for hierarchical multi-label classification. In *International Joint Conferences on Artificial Intelligence*, pages 2962–2968, 2018.
- [35] Yi Yang, Zhigang Ma, Feiping Nie, Xiaojun Chang, and Alexander G Hauptmann. Multi-class active learning by uncertainty sampling with diversity maximization. *International Journal of Computer Vision*, 113(2):113–127, 2015.
- [36] Xueying Zhan, Huan Liu, Qing Li, and Antoni B. Chan. A comparative survey: Benchmarking for pool-based active learning. In *International Joint Conferences on Artificial Intelligence*, pages 4679–4686, 2021.
- [37] Yilun Zhou, Adithya Renduchintala, Xian Li, Sida Wang, Yashar Mehdad, and Asish Ghoshal. Towards understanding the behaviors of optimal deep active learning algorithms. In *International Conference on Artificial Intelligence and Statistics*, pages 1486–1494. PMLR, 2021.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#)
 - (b) Did you describe the limitations of your work? [\[Yes\]](#) In Page 6
 - (c) Did you discuss any potential negative societal impacts of your work? [\[N/A\]](#)
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#)
 - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) Some of them are deferred to the appendix.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#)
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#)

- 442 (c) Did you report error bars (e.g., with respect to the random seed after running experiments
443 multiple times)? [Yes]
444 (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs,
445 internal cluster, or cloud provider)? [Yes] cf. appendix
446 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
447 (a) If your work uses existing assets, did you cite the creators? [Yes]
448 (b) Did you mention the license of the assets? [Yes]
449 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
450 (d) Did you discuss whether and how consent was obtained from people whose data you're us-
451 ing/curating? [Yes]
452 (e) Did you discuss whether the data you are using/curating contains personally identifiable informa-
453 tion or offensive content? [Yes]
454 5. If you used crowdsourcing or conducted research with human subjects... [No]