DIRECT ACQUISITION OPTIMIZATION FOR LOW-BUDGET ACTIVE LEARNING

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023 024

025

Paper under double-blind review

Abstract

Active Learning (AL) has gained prominence in integrating data-intensive machine learning (ML) models into domains with limited labeled data. However, its effectiveness diminishes significantly when the labeling budget is especially low. In this paper, we empirically verify the performance degradation of existing AL algorithms in the extremely low-budget settings, and then introduce *Direct Acquisition Optimization* (DAO), a novel AL algorithm that optimizes sample selections based on expected true loss reduction. Specifically, DAO utilizes influence functions to update model parameters and incorporates an additional acquisition strategy to mitigate bias in loss estimation. This approach facilitates a more accurate estimation of the overall error reduction, without extensive computations or reliance on labeled data. Experiments demonstrate the effectiveness of DAO in both low and higher budget settings, outperforming state-of-the-arts approaches across seven benchmarks.

1 INTRODUCTION

026 Active learning (AL) explores how adaptive data collection can reduce the amount of data needed by machine learning (ML) models (Settles, 2009; Schröder & Niekler, 2020; Ren et al., 2021; Zhan 027 et al., 2022). It is particularly useful when labeled data is scarce or expensive to obtain, which 028 significantly limits the adaptability of modern deep learning (DL) models due to their data-hungry 029 nature (van der Ploeg et al., 2014). In these cases, AL algorithms selectively choose the most beneficial data points for labeling, thereby maximizing the effectiveness of the training process even if the 031 data is limited in number. In fact, AL has been broadly applied in many fields (Adadi, 2021), such as medical image analysis (Budd et al., 2021), astronomy (Škoda et al., 2020), and physics (Ding et al., 033 2023), where unlabeled samples are plentiful but the process of labeling through human expert anno-034 tations or experiments is highly cost-intensive. In these contexts, judiciously selecting samples for labeling can significantly lower the expenses involved in compiling the datasets (Ren et al., 2021).

Many active learning algorithms have emerged over the past decades, with early seminal contribu-037 tions from Lewis (1995); Tong & Koller (2001); Roy & McCallum (2001), and a shift that focuses more on deep active learning – a branch of AL that targets more towards DL models in more recent years (Huang, 2021). Depending on the optimization objective, AL algorithms can be classified into 040 two categories. The first category includes heuristic objectives that are not exactly the same as the 041 evaluation metric, i.e. error reduction. Examples in this category are diversity (Sener & Savarese, 042 2017), uncertainty (Gal et al., 2017), and hybrids of both (Ash et al., 2019). Second category in-043 cludes criteria that is exactly the same as the evaluation metric, where notable approaches include 044 expected error reduction (EER) (Roy & McCallum, 2001) and its more recent follow-up works (Kil-045 lamsetty et al., 2021; Mussmann et al., 2022).

Despite the popularity of the first type of AL algorithms, existing works (Mittal et al., 2019; Hacohen et al., 2022) as well as our empirical analysis in Appendix A show that these methods often suffer heavily in low-budget settings, where the total (accumulative) sampling quota is less than 1% of the number of unlabeled data points, making them less suitable for the extreme data scarcity scenarios. In terms of the methods from the second category, their higher running time and reliance on the availability of a *validation* or *hold-out* set remain significant limitations, constraining their applicability in many data-scarcity scenarios as well. For example, EER (Roy & McCallum, 2001) re-trains the classifier for each candidate with all its possible labels, where in each time also evaluates the updated model on all the unlabeled data, making its runtime intractable especially for deep neural networks. And GLISTER (Killamsetty et al., 2021), despite being much more computationally efficient, requires a *labeled*, *hold-out* set for its sample selection process, formulated as a mixed
discrete-continuous bi-level optimization problem, to be optimized properly. While these constraints
might not be a huge limitation a few years ago, it poses a more important challenge currently as we
are adopting deep learning models to more areas, where labeled data may be extremely expensive
to acquire. More importantly, it is also worth noticing that under these scenarios, the highly limited
labeled data should have been better utilized for training than being reserved for AL algorithms.

061 Above limitations highlight a critical gap between the capabilities of current AL methodologies 062 and the urgent demands from real-world applications, underscoring the need for developing novel 063 AL strategies that can operate both relatively efficient while presenting little to none reliance on 064 the labeled set. To this end, we introduce Direct Acquisition Optimization (DAO), a novel AL algorithm that selects new samples for labeling by efficiently estimating the expected loss reduc-065 tion. Compared to EER and GLISTER, DAO solves the pain points of prohibitive running time 066 and the reliance on a separate labeled set through utilizing influence function (Ling, 1984) in model 067 parameters updates, and a more accurate, efficient unbiased estimator of loss reduction through 068 importance-weighted sampling. In summary, DAO optimizes sample selections based on expected 069 error reduction while operating efficiently through influence function-based model parameters ap-070 proximation and true overall reduced error estimation. Thorough experiments demonstrating DAO's 071 superior performance in the low-budget settings, out-performing current popular AL methods across 072 seven benchmarks. 073

074 2 RELATED WORK

075 Active learning. AL has gained a lot of attraction in recent years, with its goal to achieve bet-076 ter model performance with fewer training data (Settles, 2009; Schröder & Niekler, 2020; Ren 077 et al., 2021). There have been different selection criteria including uncertainty, diversity, query-bycommittee, version space and information-theoretic heuristics (Liu et al., 2022; Zhan et al., 2022). 079 The uncertainty-based approaches are arguably the most popular and easiest to implement, which 080 includes selection criteria such as least confidence (Lewis, 1995), minimum margin (Scheffer et al., 2001; Roth & Small, 2006; Citovsky et al., 2021), maximum entropy (Joshi et al., 2009; Settles, 081 2009) and others (Gal et al., 2017). At their core, these methods select points where the classifier is least certain. However, uncertainty-based methods can be biased towards the current learner. 083 Diversity-based methods (Settles, 2009; Bilgic & Getoor, 2009; Guo, 2010; Luo et al., 2013; Elham-084 ifar et al., 2013; Mac Aodha et al., 2014; Yang et al., 2015; Sener & Savarese, 2017; Sinha et al., 085 2019; Agarwal et al., 2020; Wu et al., 2021), on the other hand, aim to select the most representa-086 tive samples of the dataset. In addition, query-by-committee (Seung et al., 1992; Abe, 1998) and 087 version space-based (Mitchell, 1982) methods, keep a pool of models, and then select samples that 880 maximize the disagreements between them. Information-theoretic methods (Hoi et al., 2006; Barz 089 et al., 2018) typically utilize mutual information as the criterion. Hybrid method that combines both 090 uncertainty and diversity criteria, such as BADGE (Ash et al., 2019), has also been developed to 091 take advantage of both worlds. As shown later in the paper, we visually observe that the selections of our proposed DAO, although not explicitly optimized towards any of these heuristics, display 092 characteristics of an hybrid approach. 093

094 EER-based acquisition criterion. Alternatively, EER was proposed to select new training examples that result in the lowest expected error on future test examples, which directly optimizes the 096 metric by which the model will be evaluated (Roy & McCallum, 2001). In essence, EER employs 097 sample selection based on the estimated impact of adding a new data point to the training set, rather than evaluating performance against a separate validation set, meaning that it does not inherently 098 require a validation hold-out set. However, its necessity to retrain the model for every possible candidate sample and every possible label renders its cost intractable in the context of deep neural 100 networks (Budd et al., 2021; Škoda et al., 2020; Ding et al., 2023). More recent look-ahead EER-101 based AL algorithms (Mussmann et al., 2022) focus on addressing this efficiency concern. However, 102 these methods either rely on a small set of validation data to be used for the evaluation of the ex-103 pected loss reduction (Killamsetty et al., 2021), or can still be quite slow when the size of labeled 104 and unlabeled sets are large (Mohamadi et al., 2022). In this paper, we present DAO, a novel AL 105 algorithm that improves upon EER through optimizations on both model updates as well as loss 106 estimation, efficiently and effectively broadening the applicability of EER-based algorithm.

108 3 METHODOLOGY

Different from the heuristics-based AL algorithms that optimize criteria such as diversity 110 or uncertainty, DAO is built upon the EER formulation with the selection objective be-111 ing the largest reduced error evaluated on the entire unlabeled set. More specifically, 112 DAO majorly improves upon two aspects: (1) instead of re-training the classifier, we em-113 ploy influence function (Cook & Weisberg, 1982), a concept with rich history in statis-114 tical learning, to formulate the new candidate sample as a small perturbation to the ex-115 isting labeled set, so that the model parameters can be estimated without re-training; 116

and (2) instead of reserving a sepa-117 rate, relatively large labeled set for 118 validation (Killamsetty et al., 2021), 119 we sample a very small subset directly from the unlabeled set and es-120 timate the loss reduction through bias 121 correction. 122

123 Essentially, when considering each 124 candidate from the unlabeled set, we optimize the EER framework on two 125



Figure 1: Schematic of the algorithmic framework of DAO.

of its core components, which are model parameter update and true loss estimation. Additionally, we upgrade EER, which only sup-127 ports single sequential acquisition, to offer DAO in both single and batch acquisition variants by 128 incorporating stochastic samplings to the sorted estimated loss reductions. We illustrate our al-129 gorithmic framework in Fig. 1. In the following parts of this section, we first introduce a formal 130 problem statement in §3.1, and then dive into each specific component of DAO from §3.2 to §3.5. 131

3.1 PROBLEM STATEMENT 132

126

136 137

133 The optimal sequential active learning acquisition function can be formulated as selecting a budget 134 number of samples $\mathbf{x}_t^{\text{train}}$ from the current unlabeled set \mathcal{U}_t at each round t such that 135

$$\mathbf{x}_{t}^{\text{train}} = \operatorname*{arg\,min}_{\mathbf{x}_{\mathcal{S}_{i}} \subset \mathcal{U}_{t-1}} \mathbb{E}_{(y_{\mathcal{S}_{i}}|f^{*}, \mathbf{x}_{\mathcal{S}_{i}})} \left[L_{\text{true}}(f_{t|\mathbf{x}_{\mathcal{S}_{i}}, y_{\mathcal{S}_{i}}}) \right]$$
(1)

where f^* represents an optimal oracle that maps from any subset of the unlabeled data $\mathbf{x}_{S_i} \subset \mathcal{U}_{t-1}$ 138 to their ground-truth labels $y_{\mathcal{S}_i}$, and $f_{t|\mathbf{x}_{\mathcal{S}_i}, y_{\mathcal{S}_i}}$ is the model that has been trained on the union of the current labeled set \mathcal{L}_{t-1} and the current unlabeled candidates $\mathbf{x}_{\mathcal{S}_i} \subset \mathcal{U}_{t-1}$. In addition, $L_{\text{true}}(f_{t|\mathbf{x}_{\mathcal{S}_i}, y_{\mathcal{S}_i}}) = \frac{1}{|\mathcal{U}_{t-1,i}|} \sum_{\mathbf{x} \in \mathcal{U}_{t-1,i}} \ell(\mathbf{x}; f_{t|\mathbf{x}_{\mathcal{S}_i}, y_{\mathcal{S}_i}})$ represents the loss estimator that can predict the unbiased error of $f_{t|\mathbf{x}_{\mathcal{S}_i}, y_{\mathcal{S}_i}}$, where ℓ denotes the loss function. It is numerically the same as if 139 140 141 142 $f_{t|\mathbf{x}_{S_i},y_{S_i}}$ has been tested on the entire unlabeled set $\mathcal{U}_{t-1,i}$, where $\mathcal{U}_{t-1,i} = \mathcal{U}_{t-1} \setminus \mathbf{x}_{S_i}$. Such for-143 mulation represents the optimal AL criterion and aligns with any existing sequential active learning 144 algorithm — of which the goal is to select the new data points that can most significantly improve 145 the current model performance (Roy & McCallum, 2001). 146

147 Unfortunately, Eq. (1) cannot be directly implemented in practice. Because, first, we do not have 148 access to the optimal oracle f^* to reveal the labels y_{S_i} of $\mathbf{x}_{S_i} \subset \mathcal{U}_{t-1}$; second, even if we had f^* and therefore y_{S_i} , we cannot afford the cost of retraining model f_{t-1} on each $\mathcal{L}_{t-1} \cup \mathbf{x}_{S_i}$ to obtain the 149 updated $f_{t|\mathbf{x}_{S_i}, y_{S_i}}$; and third, we do not have the unbiased true loss estimator L_{true} , which demands 150 evaluating $f_{t|\mathbf{x}_{S_i}, y_{S_i}}$ on the entire $\mathcal{U}_{t-1, i}$. 151

152 Therefore, the goal of DAO is to solve the above challenges and efficiently and accurately approx-153 imate Eq. (1) for the sample selection strategy. It is also worth noting that, when $\mathbf{x}_t^{\text{train}}$ represents a 154 set of newly acquired data points, the above formulation becomes eligible for batch active learning, 155 which is more suitable for deep neural networks (Huang, 2021). 156

3.2 LABEL APPROXIMATION VIA SURROGATE 157

158 In this section, we address the first challenge when approximating Eq. (1). As we do not know the true label or true label distribution $p(y|\mathbf{x}, f^*)$ of each unlabeled sample \mathbf{x} , the best we can do is 159 provide an approximation for $p(y|\mathbf{x})$. To this end, we introduce the concept of a *surrogate* (Kossen et al., 2021), which is a model parameterized by some potentially infinite set of parameters θ . Specif-161 ically, $p(y|\mathbf{x})$ can be approximated using the marginal distribution $\pi(y|\mathbf{x}) = \mathbb{E}_{\pi(\theta)}[\pi(y|\mathbf{x},\theta)]$ with

some proposal distribution $\pi(\theta)$ over model parameters θ . In other words, we have:

$$p(y|\mathbf{x}) \approx \int_{\theta} \pi(\theta) \pi(y|\mathbf{x}, \theta) \,\mathrm{d}\theta$$
 (2)

164

183

190 191

192 193

194

195 196 197

200 201 202

205

As the sample selection process continues, new labeled points should also be used to train and update the surrogate model $\pi(\theta)$ for better approximation of the true outcomes.

168 Although ideally, a more capable surrogate is preferred for better ground truth approximations, we 169 acknowledge that the choice of surrogate model can be very sensitive to the computational con-170 straints. Therefore, if running time is at center of the concerns during sample acquisitions, using f_t 171 at step t also as the surrogate could be an efficient alternative, as we don't need to update a second model, nor do we need to run forward pass on the both models. However, this will come with the 172 cost that π_t never disagrees with f_t , which causes performance degradation for the unbiased true 173 loss estimation, which will be illustrated with more details in §3.4. Therefore, in short, we do not 174 recommend replicating f_t as surrogate in practice, unless the computational constraint is substantial. 175

176 3.3 MODEL PARAMETERS UPDATE WITHOUT RE-TRAINING

At acquisition round t, suppose we have labeled set \mathcal{L}_{t-1} and unlabeled set \mathcal{U}_{t-1} as the results from the previous round t-1, and new sample $\mathbf{x}_i \in \mathcal{U}_{t-1}$ that is currently under consideration for acquisition, the goal of this section is to estimate the parameters of model $f_{t|\mathbf{x}_i,y_i}$ that could has been obtained after training f_{t-1} on the combined dataset $\{\mathcal{L}_{t-1} \cup \mathbf{x}_i\}$. In other words, if we suppose the conventional full training converges to parameters $\hat{\theta}_{\mathbf{x}_i}$, we have:

$$\hat{\theta}_{\mathbf{x}_{i}} = \arg\min_{\theta\in\Theta} \frac{1}{|\mathcal{L}_{t-1}| + 1} \sum_{\mathbf{x}\in\{\mathcal{L}_{t-1}\cup\mathbf{x}_{i}\}} \ell(\mathbf{x};\theta)$$
(3)

where recall that $\ell(\mathbf{x}; \theta)$ denotes the loss of θ on \mathbf{x} . The core of our approach is that, instead of retraining as showed in Eq. (3), we can approximate the effect of adding a new sample as upweighting the influence function by $\frac{1}{|\mathcal{L}_{t-1}|+1}$ (Koh & Liang, 2017) and then directly estimate the updated model parameters.

Following Cook & Weisberg (1982), we have the influence function defined as:

$$\mathcal{I}_{\text{up,params}}(\mathbf{x}_i) \coloneqq \frac{d\hat{\theta}_{\epsilon,\mathbf{x}_i}}{d\epsilon} \bigg|_{\epsilon=0} = -H_{\hat{\theta}}^{-1} \nabla_{\theta} \ell(\mathbf{x}_i; \hat{\theta})$$

where $H_{\hat{\theta}}$ is the positive definite Hessian matrix (Koh & Liang, 2017). Next, we can estimate the model parameters after adding this new sample \mathbf{x}_i , as:

$$\hat{\theta}_{\mathbf{x}_{i}} - \hat{\theta} \approx \frac{1}{|\mathcal{L}_{t-1}| + 1} \mathcal{I}_{\text{up,params}}(\mathbf{x}_{i}) = -\frac{1}{|\mathcal{L}_{t-1}| + 1} H_{\hat{\theta}}^{-1} \nabla_{\theta} \ell(\mathbf{x}_{i}; \hat{\theta})$$

where $\nabla_{\theta} \ell(\mathbf{x}_i; \hat{\theta})$ could be approximated as the expected gradient of sample \mathbf{x}_i : By a slight abuse of notation of the training loss function ℓ , we denote

$$\nabla_{\theta} \ell(\mathbf{x}_{i}; \hat{\theta}) \approx \sum_{k=1}^{K} \nabla_{\theta} \ell(\mathbf{x}_{i}, \hat{y}_{k}; \hat{\theta}) \cdot \hat{p}_{k}$$
(4)

In Eq. (4), \hat{y}_k and \hat{p}_k represent model's label prediction and likelihood (e.g. confidence) respectively while K represents the total number of classes in the ground truths.

In practice, the inverse of $H_{\hat{\theta}}$ cannot be computed due to its prohibitive $O(np^2 + p^3)$ runtime (Liu et al., 2021), with *p* being the number of model parameters. The computation unavoidably becomes especially intensive when *f* is a deep neural network model (Fu et al., 2018). Luckily, we have two optimization methods, conjugate gradients (CG) (Martens et al., 2010) and stochastic estimation (Agarwal et al., 2017) at our disposal.

Conjugate gradients. As mentioned earlier, by assumption we have $H_{\hat{\theta}} \succ 0$ and $\nabla_{\theta} \ell(\mathbf{x}'; \hat{\theta})$ as a vector. Therefore, we can calculate the inverse Hessian vector product (IHVP) through first transforming the matrix inverse into an optimization problem, i.e.

214
$$H_{\hat{\theta}}^{-1} \nabla_{\theta} \ell(\mathbf{x}_{i}; \hat{\theta}) \equiv \arg \min_{t} t^{T} H_{\hat{\theta}} t - v^{T} t$$
215

and then solving it with CG (Martens et al., 2010), which speeds up the runtime effectively to O(np).

Stochastic estimation. Besides CG, we can also efficiently compute the IHVP using the stochastic estimation algorithm developed by Agarwal et al. (2017). From Neumann series, we have $A^{-1} \approx \sum_{i=0}^{\infty} (I - A)^i$ for any matrix A. Similarly, suppose we define the first j terms in the Taylor expansion of $H_{\hat{\theta}}^{-1}$ as

$$H_{\hat{\theta},j}^{-1} = \sum_{i=0}^{j} (I - H_{\hat{\theta}})^i = I + (I - H_{\hat{\theta}}) H_{\hat{\theta},j-1}^{-1}$$

we have $H_{\hat{\theta},j}^{-1} \to H_{\hat{\theta}}^{-1}$ as $j \to \infty$. The core idea of the stochastic estimation is that the Hessian matrix $H_{\hat{\theta}}$ can be substituted with any unbiased estimation when computing $H_{\hat{\theta}}^{-1}$. In practice, we sample n_{ihvp} data points from the existing labeled set \mathcal{L}_{t-1} and use $\nabla_{\theta}^2 \ell(\mathbf{x}_i; \hat{\theta})$ as the estimator of $H_{\hat{\theta}}$ (Liu et al., 2021). Notice that since n_{ihvp} is usually very small (in our experiments we used $n_{\text{ihvp}} = 8$), it does not create a constraint on the size of the current labeled set, which does not interfere with the low-budget settings.

Finally, we can approximate the model parameters after the addition of x_i as

$$\hat{\theta}_{\mathbf{x}_{i}} = \hat{\theta} - \frac{1}{n+1} H_{\hat{\theta}}^{-1} \nabla_{\theta} \ell(\mathbf{x}_{i}; \hat{\theta})$$
(5)

which does not require any re-training. And we will demonstrate in §5.2 that this parameter update strategy provides much better approximations than the naive single backpropagation as seen in the existing AL literature (Killamsetty et al., 2021).

2372383.4 EFFICIENT UNBIASED LOSS ESTIMATION

Referring back to Eq. (1), the last challenge that we need to address is to gain access to the unbiased true loss estimator L_{true} . In other words, we want to predict the *true* performance of $f_{t|\mathbf{x}_{i},y_{i}}$ on the unlabeled set $\mathcal{U}_{t,i}$ without exhaustive testing. Strictly, such evaluation cannot be drawn until $f_{t|\mathbf{x}_{i},y_{i}}$ is evaluated on the entire unlabeled set $\mathcal{U}_{t,i}$. However, this is infeasible in practice.

Such approximation is typically carried out in other approaches (Killamsetty et al., 2021; Mussmann et al., 2022) by randomly sampling a labeled validation set \mathcal{V} at the beginning of the entire acquisition process, which will later be used for evaluations in all the subsequent acquisition episodes. Despite the simplicity as well as being i.i.d., which makes the estimated loss unbiased by nature, this approximation method suffers from large variance as the size of \mathcal{V} is usually much smaller than \mathcal{U} , which unavoidably hurts the acquisition performance. It is also contradictory to the goal of AL in general, especially under the low-budget settings, as discussed in §1.

Different from the existing works, we propose to sample a subset C from current \mathcal{L}_{t-1} in each acquisition round based on an alternative acquisition function, and then correct the bias in the loss induced from this acquisition function. In the meantime, we also want to keep the variance low, so that the final corrected loss enjoys both low bias and low variance, which is more preferable than the zero bias but high variance that the random i.i.d. sampling has.

Specifically, continuing with the notations from §3.1, let $C = {\mathbf{x}_{t,1}, \ldots, \mathbf{x}_{t,m}, \ldots, \mathbf{x}_{t,n_c}}$, where $C \subset U_{t-1}$, be the subset containing n_C samples selected for this true loss estimation at each round t. Farquhar et al. (2021) shows that if $\mathbf{x}_{t,m}$ is sampled in proportion to the true loss of each data point, the bias originated from this selection can be corrected through the Monte Carlo estimator \hat{R}_{LURE}^{-1} . Following our notations, it takes the form:

265 266 267

269

221 222

223 224

225

226

227

228

229

230

232 233 234

235

236

$$\hat{R}_{\text{LURE}} = \frac{1}{n_{\mathcal{C}}} \sum_{m=1}^{n_{\mathcal{C}}} v_m \ell\left(\mathbf{x}_{t,m}; f\right)$$
(6)

where recall that ℓ denotes the loss of f, and the importance weight v_m is

$$v_m = 1 + \frac{|\mathcal{U}_{t-1}| - n_{\mathcal{C}}}{|\mathcal{U}_{t-1}| - m} \left(\frac{1}{(|\mathcal{U}_{t-1}| - m + 1)q_t^*(m)} - 1 \right)$$
(7)

with $q_t^*(m)$ being the acquisition distribution of index m at round t. Importantly, the variance can be significantly reduced if the acquisition distribution $q_t^*(m)$ is proportion to the true loss of each

¹LURE stands for Levelled Unbiased Risk Estimator

273

274

279

284

285 286

287 288

289

295

270 data point. Again, this is not feasible as we do not have access to the labels for \mathcal{U}_{t-1} . However, 271 following Kossen et al. (2021), we can approximate $q_t^*(m)$ with 272

$$q_t(m) = -\sum_y \pi(y|\mathbf{x}_{t,m}) \log f(\mathbf{x}_{t,m})$$

275 for classification tasks when the loss function is the cross-entropy loss, and where π is conveniently 276 just our surrogate discussed in §3.2. Referring back to the discussion we had on choosing a good 277 surrogate π , with $f(\mathbf{x})$ being designed to approximate $p(y|\mathbf{x})$ as well, the surrogate π should ideally 278 be different from f so that more diversity is introduced in the acquisitions.

To put all components together, our loss correction process involves selecting samples in C following

$$\mathbf{x}_{t,m} \propto -\sum_{y} \pi_{t-1}(y|\mathbf{x}) \log f_{t-1}(y|\mathbf{x})$$
(8)

where π_{t-1} is the surrogate model π at round t-1. Finally, the corrected loss s_i can be approximated using \hat{R}_{LURE} as

$$s_i = \frac{1}{n_{\mathcal{C}}} \sum_{m=1}^{n_{\mathcal{C}}} \hat{v}_m \ell\left(\mathbf{x}_{t,m}; f_t\right)$$
(9)

where \hat{v}_m , which depends on the choice of $\mathbf{x}_{t,m}$, is the approximated version of the original v_m defined in Eq. (7). Specifically, \hat{v}_m takes the form

$$\hat{v}_m = 1 + \frac{|\mathcal{U}_{t-1}| - n_{\mathcal{C}}}{|\mathcal{U}_{t-1}| - m} \left(\frac{1}{(|\mathcal{U}_{t-1}| - m + 1)q_t(m)} - 1 \right)$$
(10)

where $q_t(m)$ is the acquisition function defined in Eq. (8).

3.5 BATCH ACQUISITION VIA STOCHASTIC SAMPLING

296 In §3.1, we briefly discussed that when $\mathbf{x}_t^{\text{train}}$ represents a set of data points (instead of a 297 single one), the formulation in Eq. (1) essentially represents the batch active learning sce-298 nario. Suppose the acquisition budget per round is k, although selecting the top k samples 299 with the lowest estimated losses (or highest expected error reduction) is straightforward, this approach is sub-optimal. This is because top-k acquisition, while effective to some degree 300 due to its greedy nature, overlooks the crucial interactions among data points in batch acqui-301 sitions. Specifically, while aiming to select the most informative unlabeled points, top-k ac-302 quisition may lead to redundant choices, diminishing the overall benefit of the acquisition. 303

Inspired by Kirsch et al. (2021), we pro-304 pose to similarly perturb the original 305 ranking of the estimated true losses so 306 that the batch sampling provides better 307 acquisitions when the most informative 308 data points may be duplicated. Sup-309 pose at acquisition episode t, we rank 310 the set of estimated true loss of each unlabeled data point in ascending orders 311 as $\{l_{ ext{true},i}\}_{\mathbf{x}_i\in\mathcal{U}_{t-1}}, ext{ such that } l_{ ext{true},i} \leq$ 312 313 $\hat{l}_{\text{true.}i}, \forall i \leq j \text{ and } \mathbf{x}_i, \mathbf{x}_j \in \mathcal{U}_{t-1}, \text{ we can}$ 314 perturb the ranking with three strategies: 315 soft-rank, soft-max, and power acquisition, to improve batch performance from 316 the naive top-k sampling. 317

318 Soft-rank acquisition. Soft-rank ac-319 quisition relies on the relative ordering 320 of the scores while ignoring the abso-321 lute score values. It samples the data point ranked at index i with probabil-322 ity $p_{\text{softrank}}(i) = i^{-\beta}$, where β is the 323 "coldness" parameter and is kept as 1

Algorithm 1 Direct	Acqu	isitio	n Opti	mization (DAO)
input Episode t, un	label	ed se	t \mathcal{U}_{t-1}	, labeled	set \mathcal{L}_{t-1} ,
model f_{t-1} , sur	rogat	e π_{t-}	$_1$, bud	lget $k, n_{\rm ih}$	vp (§3.3),
output Acquisition	set	\mathcal{A}_t	=	$\{\mathbf{x}_{t,1}^{\text{train}},.\}$	$\ldots, \mathbf{x}_{t,k}^{\text{train}} \}$

 \triangleright Eq. (1) Approximate $p(u|\mathbf{x})$ for all $\mathbf{x} \in \mathcal{U}_{t-1} \triangleright$ §3.2. Eq. (2)

1: Approximate
$$p(y|\mathbf{x})$$
 for all $\mathbf{x} \in \mathcal{U}_{t-1} \triangleright$ §3.2, Eq. (2
2: Initialize array S where $|S| = |\mathcal{U}_{t-1}|$

2: Initialize array S where
$$|S|$$

- 3: for i = 1 to $|\mathcal{U}_{t-1}|$ do 4:
 - Let $\mathcal{U}_{t,i} = \mathcal{U}_{t-1} \setminus {\mathbf{x}_i}$
- Randomly sample n_{ihvp} data points from $U_{t,i}$ 5: Approximate parameters of $f_{t|\mathbf{x}_i, y_i}$ 6: ⊳ §3.3,
- Eq. (5) 7: Acquire n_c samples from $\mathcal{U}_{t,i}$
 - ▷ §3.4, Eq. (8) Compute s_i and add to S ⊳ §3.4, Eq. (9)

9: end for

8:

- 10: Sort S in ascending order
- 11: **if** *k* > 1 **then**
- 12: Perturb S ▷ Methods showed in §3.5 13: end if
- 14: Return top-k samples in S as A_t

throughout this paper. It is not hard to notice that $p_{\text{softrank}}(i)$ is invariant to $\hat{l}_{\text{true},i}$, as long as the relative ranking remains the same. More conveniently, with sampled Gumbel noise $\epsilon_i \sim$ Gumbel(0; β^{-1}), taking the top-k data points from the perturbed ranked list

$$\hat{l}_{\text{true},i}^{\text{softrank}} = -\log i + \epsilon_i$$

is equivalent to sampling $p_{\text{softrank}}(i)$ without replacement (Huijben et al., 2022).

Soft-max acquisition. In contrast to soft-rank, soft-max acquisition uses the actual scores, i.e., the
 estimated true losses, instead of their relative orderings. However, this acquisition does not rely on
 the semantics of the actual values, resulting in the transformed true loss simply being:

$$\hat{l}_{\text{true},i}^{\text{softmax}} = \hat{l}_{\text{true},i} + \epsilon_i$$

where ϵ_i remains the same Gumbel noise as in the soft-rank acquisition. Statistically, choosing the top-k data points from this perturbed ranked list is equivalent to sample from $p_{\text{softmax}}(i) = e^{\beta i}$ without replacement.

Power acquisition. While neither soft-rank or soft-max acquisitions take the semantic meaning of the actual score values into account when designing the acquisition distribution, power acquisition uses the value directly when determining the perturbed values. Specifically, the power acquisition perturbs the scores as

$$\hat{l}_{\text{true},i}^{\text{power}} = \log \hat{l}_{\text{true},i} + \epsilon_i$$

where again ϵ_i is the Gumbel noise, and choosing the top-k indices from this new list is equivalent to sampling from $p_{\text{power}}(i) = i^{\beta}$ without replacement. Results comparing DAO with different batch acquisition strategies discussed above are showed in Appendix B. Combining all the components, the pseudocode of DAO is summarized in Algorithm 1.

347 4 EXPERIMENTS

328

333

We evaluate DAO on seven classification benchmarks including digit recognition datasets
MNIST (LeCun et al., 1998), Street-View House Numbers recognition (SVHN) (Sermanet et al., 2012), object classification datasets STL-10 (Coates et al., 2011), CIFAR-10, CIFAR-100 (Krizhevsky et al., 2009), as well as domain-specific datasets Fashion-MNIST (Xiao et al., 2017) and Stanford Cars (Cars196) (Krause et al., 2013).

354 4.1 EXPERIMENTAL SETUP

355 **Baselines.** To ensure fair comparisons, besides baseline methods that we empirically surveyed in 356 Appendix A, we also include other state-of-the-arts AL methods, including Deep Bayesian Active 357 Learning (DBAL) (Gal et al., 2017) and GLISTER (Killamsetty et al., 2021), where GLISTER is a direct competitor that also optimizes the EER framework. For all the baselines, we used the 358 default/recommended parameters and their official implementations if publically available. In terms 359 of earlier works such as least confidence (Lewis, 1995), minimum margin (Scheffer et al., 2001), 360 and maximum entropy (Settles, 2009), we used the peer-reviewed deep active learning framework 361 DeepAL+ (Zhan et al., 2022). All experiments are repeated ten times with different random seeds. 362

Implementation details. Throughout the section, we set ResNet-18 (He et al., 2016) as the model f to be trained from scratch. We employed VGG16 (Simonyan & Zisserman, 2014), initialized with random weights, as our surrogate π . We used stochastic estimation (Agarwal et al., 2017) when estimating the updated model parameters, as discussed in §3.3. We choose $n_{\text{ihvp}} = 8$ when approximating the unbiased estimator of $H_{\hat{\theta}}$, and set $n_{\mathcal{C}} = 16$ for biased loss correction as in §3.4.

368 4.2 Performance under Low Budgets

369 Digit recognition. First, we demonstrate DAO's effectiveness through two digit recognition bench-370 marks: MNIST (LeCun et al., 1998) and SVHM (Sermanet et al., 2012). MNIST is a collection of 371 handwritten digits consisting of 60k training and 10k test images, while SVHN is a more challeng-372 ing dataset containing over 600k real-world house numbers images taken from street views. Both 373 datasets contain 10 classes corresponding to digits from 0 to 9. Based on insights from Appendix A, 374 we define a general rule of low-budget setting as one image per class, which translates to initial label size $|\mathcal{L}_{\text{init}}^{\text{MNIST}}| = 10$ and per-episode budget $B_{\text{MNIST}} = 10$ for MNIST. Given that SVHN is 375 more challenging, and there are ten times more unlabeled images than in MNIST (600k vs. 60k), we 376 experiment both 10 and 100 for our initial labeled size and per-round budget for SVHN. The results 377 are shown in Fig. 2b and 2c.

378

379380381382

384 385

386

387

394

396

397

398



Figure 2: Experiment results comparing DAO with existing AL algorithms across seven benchmarks. In all subplots, horizontal axis represents the accumulative size of the labeled set, while vertical axis indicates classification accuracy.

399 **Object classification.** Next, we assess DAO on more general and complex object classification 400 tasks. STL-10 (Coates et al., 2011) is a benchmark dataset derived from labeled examples in the 401 ImageNet (Deng et al., 2009). Specifically, STL-10 contains 5k labeled 96×96 color images spread 402 across 10 classes, as well as 8k images in the test split. CIFAR-10 (Krizhevsky et al., 2009) contains 403 a collection of 60k 32x32 color images in 10 different classes, with 6k images per class. CIFAR-100 404 is similar to CIFAR-10, but covers a much wider range, containing 100 classes where each class holds 600 images. Continuing with the low-budget setting (1 image per class), we have $|\mathcal{L}_{init}^{STL-10}| =$ 405 10, $B_{\text{STL-10}} = 10$ for STL-10, $|\mathcal{L}_{\text{init}}^{\text{CIFAR-10}}| = 10$, $B_{\text{CIFAR-10}} = 10$ for CIFAR-10 and $|\mathcal{L}_{\text{init}}^{\text{CIFAR-100}}| = 10$ 406 100, $B_{\text{CIFAR-100}} = 100$ for CIFAR-100. The results are shown in Fig. 2d, 2e and 2f respectively. 407

408 **Domain Specific Tasks** The last part of our experiments involves case studies on applying DAO to 409 domain-specific tasks, which simulates many real-world applications. Specifically, we use Fashion-410 MNIST (Xiao et al., 2017) and StanfordCars (Krause et al., 2013), also known as Cars196, in this 411 experiment. FashionMNIST is structure-wise similar to MNIST, comprising 28×28 images of 70k 412 fashion products from 10 categories, with 7k images per category. The training set contains 60k images, while the test set includes the rest. StanfordCars is a large collection of car images, containing 413 16,185 images with a near-balanced ratio on the train/test split, resulting in 8,144 and 8,041 images 414 for training and testing. There are 196 classes in total, where each class consists of the year, make, 415 model of a car (e.g., 2012 Tesla Model S). The results of both datasets are shown in Fig. 2g and 2h. 416

417 Results discussion. From Fig. 2, we notice that the proposed DAO outperforms popular 418 AL state-of-the-arts by a clear margin across 419 all seven benchmarks. Especially, with SVHN, 420 when the budget is extremely low (B = 10,421 which is 0.0017% of the unlabeled size), DAO 422 leads the performance by a very large gap, indi-423 cating its superior capability in the low-budget 424 setting. Such performance does not degrade 425 much as the budget constraint is relaxed. As 426 shown in Fig. 2c, DAO still performs relatively 427 well. The only experiment that DAO does not 428 improve as much is the StanfordCars. However, the accuracy improvement from DAO is more 429 smooth and has less variance, indicating better 430 robustness when applied to the more challeng-431 ing (StanfordCars has 196 classes) applications.



Figure 3: Higher-budget experiment results comparing DAO with existing AL algorithms. In both subplots, horizontal axis represents the accumulative size of the labeled set, while vertical axis indicates classification accuracy.

432 4.3 PERFORMANCE UNDER HIGHER BUDGETS

To further evaluate the capabilities of the proposed DAO beyond its targeting focus on low-budget AL, we conducted additional experiments with higher budgets on the CIFAR-100 and STL10 datasets. Specifically, we keep the same experimental settings as previously, with the same initial labeled set and per-round acquisition budget, and repeated the process five times with different random seeds. However, we extended the number of rounds in each acquisition from 10 to 50, increasing the budget by five times. And to make the plot more clear, we plot every five rounds. As shown in Fig. 3, DAO consistently outperforms all other baselines throughout the entire acquisition process, demonstrating superior performance across both low and high budgets.

441

442 443

444

5

5.1 ABLATIONS ON DIFFERENT SURROGATES

COMPONENT ANALYSIS AND ABLATION STUDIES

To further understand how the surrogate model impacts the performance of DAO, we conducted 445 additional experiments on CIFAR-10 using different surrogates. Specifically, we compared: 446 (1) VGG16: The surrogate model currently used in the paper draft. (2) ResNet18: An efficient 447 variant of DAO where the main model under training serves as the surrogate, i.e., $\pi_t = f_t$. In this 448 case, the equation for approximating $q_t^*(m)$, the acquisition distribution of candidate index m at 449 round t (the equation between lines 225 and 226), reduces from cross-entropy to entropy. (3) Sim-450 *pleCNN*: A simpler version of ResNet18 with six convolution blocks, each containing one Conv 451 layer followed by BatchNorm and ReLU. and (4) Oracle: An unrealistic setting assuming access to 452 an oracle surrogate model. We used the same experimental setup with CIFAR-10 as in our draft: 453 starting with 10 labeled samples, acquiring 10 samples per round, and continuing for 10 rounds.

The was repeated five times with different random seeds.

455 As shown in Fig. 4, DAO with oracle surrogate performs 456 the best, followed by VGG16 and SimpleCNN. ResNet18 457 performs the worst among all DAO variants, which aligns 458 with our expectation that performance degrades when π_t 459 never disagrees with f_t . However, all DAO variants out-460 perform the random and GLISTER baselines by a clear margin. It is worth noting that while the oracle surrogate 461 achieves the best results, the improvement over VGG16 462 and SimpleCNN is not substantial. We think this is likely 463 because, when selecting samples for the unbiased loss 464 estimation, the acquisition distribution $q_t^*(m)$ approxi-465 mated with $q_t(m) = -\sum_y \pi(y|\mathbf{x}_{t,m}) \log f(\mathbf{x}_{t,m})$, does 466 not solely depend on the surrogate quality. Although in 467 the unrealistic case of an oracle surrogate, this creates an



Figure 4: DAO performance when using different surrogate models.

disadvantage, but in practical scenarios, this approximation provides better robustness in prevent ing the negative impact of poor quality of surrogate on unbiased loss estimation, especially in early
 stages where models might overfit.

Table 1: Surrogate model performance on unla-beled set during acquisition.

Budget	20	40	60	80	100
SimpleCNN	13.79	13.91	15.11	15.01	14.80
VGG16	18.50	20.07	21.02	22.44	22.75
ResNet18	18.60	21.79	23.86	25.36	27.05

478 cally optimizes the sample selection that enhances the performance of the target model f.

479 480

5.2 ACCURACY ON MODEL APPROXIMATION

First, we assess if estimating the model parameters updates through modelling the effect of adding a new sample as upweighting the influence function provides a more accurate model performance approximation than using single backpropagation as seen in GLISTER (Killamsetty et al., 2021). Specifically, we conduct the experiments on CIFAR-10 (Krizhevsky et al., 2009), with initial labeled size $|\mathcal{L}_{init}^{CIFAR-10}| = 100$ (randomly sampled from the train split), per-episode budget $B_{CIFAR-10} = 1$, and number of acquisition episode E = 25. We compare the updated models performance (accu-

486 racy) on the test split of CIFAR-10. Different from the experiments in §4, we do not apply any 487 AL algorithm when acquiring the sample in each round. Instead, we randomly choose B sample 488 in each acquisition round from the unlabeled set and then update the models through both methods 489 with the same selected sample. To access the difference between models updated with our influence 490 function-based method and single backpropagation, we compute the mean squared error (MSE) between the performance of each model and the model updated by conventional full training, which is 491 defined in Eq. (3). As shown in Fig. 5a, the proposed method provides more accurate (smaller mean 492 and median) and more robust (smaller std.) model approximations than single backpropagation, 493 contributing to the performance gain we observe in §4. 494

495 5.3 BIAS CORRECTION VS. RANDOM SAMPLING

Next, we conduct ablation studies on replacing the proposed loss estimation (§3.4) with the average loss of randomly sampled data points. More specifically, we replace the estimated loss s_i from averaging the corrected loss (Eq. (6)) of the acquired samples via an alternative acquisition criteria (Eq. (8)) with averaging losses of the samples acquired uniformly, i.e., at round t, we have $s_i^{\text{random}} = \frac{1}{M_{\text{random}}} \sum_{m=1}^{M_{\text{random}}} \ell(\mathbf{x}_{t,m}; f_t)$ where $\mathbf{x}_{t,m} \sim U(1, |\mathcal{U}_{t,i}|)$. We choose two $M_{\text{random}} = 16$ and 256, where former provides a direct comparison with our proposed loss estimation approach, and latter represents a brute-force solution that works relatively well but is often infeasible in practice due to intensive running time. The results are shown in Fig. 5b.

504 We see that the proposed method 505 performs even better than the conventional random-sampling loss es-506 timation with large sampling size, 507 while computationally being only 508 1/8 of the run time. Additionally, 509 the variance of our method is much 510 smaller, indicating more robust loss 511 estimation and thus more robust ac-512 quisition performance. 513



514 5.4 DIFFERENT BATCH

515 ACQUISITION STRATEGIES

516 We conducted additional ablation
517 studies comparing various stochas518 tic sampling methods as detailed in
\$3.5. Results are documented in Ap-



pendix B. Our findings reveal that the proposed algorithm, even when simply selecting the top ksamples without applying any of the stochastic strategies, outperforms existing methods. Performance further improves with the implementation of these sampling strategies. It is important to note that we have not designed specific sampling strategies for our algorithm; instead, we utilized existing methods to showcase the efficacy of DAO framework.

525 5.5 INTERPRETING DAO WITH OTHER AL CRITERIA

In this section, we analyze the criterion optimized by DAO and compare it to common criteria such as diversity and uncertainty, using visual representations of the data samples collected by DAO. The detailed plots are available in Appendix C. Throughout multiple acquisition rounds, the data selected by DAO demonstrate notable diversity with uniform distribution across the sample space. However, in contrast to traditional uncertainty-based methods, selections within a single round by DAO also incorporate elements of uncertainty. This hybrid approach explains the performance improvements observed in §4 over algorithms that solely focus on diversity or uncertainty.

533 6 CONCLUSIONS

In this paper, we introduced Direct Acquisition Optimization, a novel algorithm designed to optimize sample selections in low-budget settings. DAO hinges on the utilization of influence functions for model parameter updates and a separate acquisition strategy to mitigate bias in loss estimation, represents a significant optimization of the EER method and its modern follow-ups. Through empirical studies, DAO has demonstrated superior performance in both low and higher budget settings, outperforming existing methods by a significant margin across seven datasets.

540 REFERENCES 541

552

553

554

555 556

558

559

561

565

566

567

576

580

581

582

- 542 Naoki Abe. Query learning strategies using boosting and bagging. In International Conference on Machine Learning, 1998, pp. 1–9, 1998. 543
- 544 Amina Adadi. A survey on data-efficient algorithms in big data era. Journal of Big Data, 8(1):24, 2021. 546
- Naman Agarwal, Brian Bullins, and Elad Hazan. Second-order stochastic optimization for machine 547 learning in linear time. The Journal of Machine Learning Research, 18(1):4148–4187, 2017. 548
- 549 Sharat Agarwal, Himanshu Arora, Saket Anand, and Chetan Arora. Contextual diversity for active 550 learning. In Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 551 23-28, 2020, Proceedings, Part XVI 16, pp. 137-153. Springer, 2020.
 - Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep batch active learning by diverse, uncertain gradient lower bounds. arXiv preprint arXiv:1906.03671, 2019.
 - Björn Barz, Christoph Käding, and Joachim Denzler. Information-theoretic active learning for content-based image retrieval. In German Conference on Pattern Recognition, pp. 650-666. Springer, 2018.
- Mustafa Bilgic and Lise Getoor. Link-based active learning. In NIPS Workshop on Analyzing 560 Networks and Learning with Graphs, volume 4, pp. 9, 2009.
- 562 Samuel Budd, Emma C Robinson, and Bernhard Kainz. A survey on active learning and human-in-563 the-loop deep learning for medical image analysis. *Medical Image Analysis*, 71:102062, 2021.
 - Akshay L Chandra, Sai Vikas Desai, Chaitanya Devaguptapu, and Vineeth N Balasubramanian. On initial pools for deep active learning. In NeurIPS 2020 Workshop on Pre-registration in Machine Learning, pp. 14-32. PMLR, 2021.
- 568 Gui Citovsky, Giulia DeSalvo, Claudio Gentile, Lazaros Karydas, Anand Rajagopalan, Afshin Ros-569 tamizadeh, and Sanjiv Kumar. Batch active learning at scale. Advances in Neural Information 570 Processing Systems, 34:11933–11944, 2021.
- 571 Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised 572 feature learning. In Proceedings of the fourteenth international conference on artificial intelli-573 gence and statistics, pp. 215–223. JMLR Workshop and Conference Proceedings, 2011. 574
- 575 R Dennis Cook and Sanford Weisberg. Residuals and influence in regression, 1982.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-577 erarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, 578 pp. 248-255. Ieee, 2009. 579
 - Yongcheng Ding, José D Martín-Guerrero, Yolanda Vives-Gilabert, and Xi Chen. Active learning in physics: From 101, to progress, and perspective. Advanced Quantum Technologies, pp. 2300208, 2023.
- 583 Ehsan Elhamifar, Guillermo Sapiro, Allen Yang, and S Shankar Sasrty. A convex optimization 584 framework for active learning. In Proceedings of the IEEE International Conference on Computer 585 Vision, pp. 209–216, 2013. 586
 - Sebastian Farquhar, Yarin Gal, and Tom Rainforth. On statistical bias in active learning: How and when to fix it. arXiv preprint arXiv:2101.11665, 2021.
- 589 Weijie Fu, Meng Wang, Shijie Hao, and Xindong Wu. Scalable active learning by approximated 590 error reduction. In Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining, pp. 1396–1405, 2018. 592
- Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In International conference on machine learning, pp. 1183–1192. PMLR, 2017.

594

595 cessing Systems, 23, 2010. 596 Guy Hacohen, Avihu Dekel, and Daphna Weinshall. Active learning on a budget: Opposite strategies 597 suit high and low budgets. arXiv preprint arXiv:2202.02794, 2022. 598 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-600 nition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 601 770-778, 2016. 602 Steven CH Hoi, Rong Jin, Jianke Zhu, and Michael R Lyu. Batch mode active learning and its 603 application to medical image classification. In Proceedings of the 23rd international conference 604 on Machine learning, pp. 417-424, 2006. 605 606 Kuan-Hao Huang. Deepal: Deep active learning in python. arXiv preprint arXiv:2111.15258, 2021. 607 608 Iris AM Huijben, Wouter Kool, Max B Paulus, and Ruud JG Van Sloun. A review of the gumbelmax trick and its extensions for discrete stochasticity in machine learning. IEEE Transactions on 609 Pattern Analysis and Machine Intelligence, 45(2):1353–1371, 2022. 610 611 Ajay J Joshi, Fatih Porikli, and Nikolaos Papanikolopoulos. Multi-class active learning for image 612 classification. In 2009 ieee conference on computer vision and pattern recognition, pp. 2372-613 2379. IEEE, 2009. 614 615 Krishnateja Killamsetty, Durga Sivasubramanian, Ganesh Ramakrishnan, and Rishabh Iyer. Glister: Generalization based data subset selection for efficient and robust learning. In Proceedings of the 616 AAAI Conference on Artificial Intelligence, volume 35, pp. 8110–8118, 2021. 617 618 Andreas Kirsch, Sebastian Farquhar, Parmida Atighehchian, Andrew Jesson, Frederic Branchaud-619 Charron, and Yarin Gal. Stochastic batch acquisition for deep active learning. arXiv preprint 620 arXiv:2106.12059, 2021. 621 622 Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In International conference on machine learning, pp. 1885–1894. PMLR, 2017. 623 624 Jannik Kossen, Sebastian Farquhar, Yarin Gal, and Tom Rainforth. Active testing: Sample-efficient 625 model evaluation. In International Conference on Machine Learning, pp. 5753–5763. PMLR, 626 2021. 627 Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained 628 categorization. In Proceedings of the IEEE international conference on computer vision work-629 shops, pp. 554-561, 2013. 630 631 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 632 2009. 633 634 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. 635 636 David D Lewis. A sequential algorithm for training text classifiers: Corrigendum and additional 637 data. In Acm Sigir Forum, volume 29, pp. 13-19. ACM New York, NY, USA, 1995. 638 639 Robert F Ling. Residuals and influence in regression, 1984. 640 Peng Liu, Lizhe Wang, Rajiv Ranjan, Guojin He, and Lei Zhao. A survey on active deep learning: 641 from model driven to data driven. ACM Computing Surveys (CSUR), 54(10s):1-34, 2022. 642 643 Zhuoming Liu, Hao Ding, Huaping Zhong, Weijia Li, Jifeng Dai, and Conghui He. Influence selec-644 tion for active learning. In Proceedings of the IEEE/CVF International Conference on Computer 645 Vision, pp. 9274–9283, 2021. 646

Yuhong Guo. Active instance sampling via matrix partition. Advances in Neural Information Pro-

647 Wenjie Luo, Alex Schwing, and Raquel Urtasun. Latent structured active learning. Advances in Neural Information Processing Systems, 26, 2013.

648 649 650	Oisin Mac Aodha, Neill DF Campbell, Jan Kautz, and Gabriel J Brostow. Hierarchical subquery evaluation for active learning on a graph. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 564–571, 2014.
651 652 653	James Martens et al. Deep learning via hessian-free optimization. In <i>ICML</i> , volume 27, pp. 735–742, 2010.
654	Tom M Mitchell. Generalization as search. Artificial intelligence, 18(2):203–226, 1982.
655	Sudharshu Mittal Marin Tatarsharla Özzür Ciash and Themas Dear Deating with illusions
656 657	about deep active learning. <i>arXiv preprint arXiv:1912.05361</i> , 2019.
658 659 660	Mohamad Amin Mohamadi, Wonho Bae, and Danica J Sutherland. Making look-ahead active learn- ing strategies feasible with neural tangent kernels. <i>Advances in Neural Information Processing</i> <i>Systems</i> , 35:12542–12553, 2022.
661 662	Stephen Mussmann, Julia Reisler, Daniel Tsai, Ehsan Mousavi, Shayne O'Brien, and Moises Gold- szmidt. Active learning with expected error reduction. <i>arXiv preprint arXiv:2211.09283</i> , 2022.
664 665 666	Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Brij B Gupta, Xiaojiang Chen, and Xin Wang. A survey of deep active learning. <i>ACM computing surveys (CSUR)</i> , 54(9):1–40, 2021.
667 668 669 670	Dan Roth and Kevin Small. Margin-based active learning for structured output spaces. In Ma- chine Learning: ECML 2006: 17th European Conference on Machine Learning Berlin, Germany, September 18-22, 2006 Proceedings 17, pp. 413–424. Springer, 2006.
671 672	Nicholas Roy and Andrew McCallum. Toward optimal active learning through monte carlo estima- tion of error reduction. <i>ICML, Williamstown</i> , 2:441–448, 2001.
673 674 675	Tobias Scheffer, Christian Decomain, and Stefan Wrobel. Active hidden markov models for infor- mation extraction. In <i>International symposium on intelligent data analysis</i> , pp. 309–318. Springer, 2001.
676 677 678	Christopher Schröder and Andreas Niekler. A survey of active learning for text classification using deep neural networks. <i>arXiv preprint arXiv:2008.07267</i> , 2020.
679 680	Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. <i>arXiv preprint arXiv:1708.00489</i> , 2017.
681 682 683 684	Pierre Sermanet, Soumith Chintala, and Yann LeCun. Convolutional neural networks applied to house numbers digit classification. In <i>Proceedings of the 21st international conference on pattern recognition (ICPR2012)</i> , pp. 3288–3291. IEEE, 2012.
685	Burr Settles. Active learning literature survey. 2009.
686 687 688	H Sebastian Seung, Manfred Opper, and Haim Sompolinsky. Query by committee. In <i>Proceedings</i> of the fifth annual workshop on Computational learning theory, pp. 287–294, 1992.
689 690	Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. <i>arXiv preprint arXiv:1409.1556</i> , 2014.
691 692 693 694	Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell. Variational adversarial active learning. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 5972–5981, 2019.
695 696	Petr Škoda, Ondřej Podsztavek, and Pavel Tvrdík. Active deep learning method for the discovery of objects of interest in large spectroscopic surveys. <i>arXiv preprint arXiv:2009.03219</i> , 2020.
697 698 699	Simon Tong and Daphne Koller. Support vector machine active learning with applications to text classification. <i>Journal of machine learning research</i> , 2(Nov):45–66, 2001.
700 701	Tjeerd van der Ploeg, Peter C Austin, and Ewout W Steyerberg. Modern modelling techniques are data hungry: a simulation study for predicting dichotomous endpoints. <i>BMC medical research methodology</i> , 14(1):1–13, 2014.

702 703 704 705	Tsung-Han Wu, Yueh-Cheng Liu, Yu-Kai Huang, Hsin-Ying Lee, Hung-Ting Su, Ping-Chia Huang, and Winston H Hsu. Redal: Region-based and diversity-aware active learning for point cloud semantic segmentation. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pp. 15510–15519, 2021.
706 707 708	Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmark- ing machine learning algorithms. <i>arXiv preprint arXiv:1708.07747</i> , 2017.
709 710 711	Yi Yang, Zhigang Ma, Feiping Nie, Xiaojun Chang, and Alexander G Hauptmann. Multi-class active learning by uncertainty sampling with diversity maximization. <i>International Journal of Computer Vision</i> , 113:113–127, 2015.
712 713 714	Xueying Zhan, Qingzhong Wang, Kuan-hao Huang, Haoyi Xiong, Dejing Dou, and Antoni B Chan. A comparative survey of deep active learning. <i>arXiv preprint arXiv:2203.13450</i> , 2022.
715 716 717	Yu Zhu, Jinghao Lin, Shibi He, Beidou Wang, Ziyu Guan, Haifeng Liu, and Deng Cai. Addressing the item cold-start problem by attribute-driven active learning. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 32(4):631–644, 2019.
718 719 720	
721 722	
723 724	
725 726 727	
728 729	
730 731 732	
733 734	
735 736	
737 738 739	
740 741	
742 743	
744 745 746	
747 748	
749 750	
751 752 753	
754 755	

756 APPENDIX

758

A LOW-BUDGET ACTIVE LEARNING: A MOTIVATING CASE STUDY

In this section, we provide an empirical analysis to demonstrate that commonly used heuristic-based AL algorithms do not work well under very low-budget settings. Specifically, we analyze (1) uncertainty sampling methods including least confidence (Lewis, 1995), minimum margin (Scheffer et al., 2001), maximum entropy (Settles, 2009), and Bayesian Active Learning by Disagreement (BALD) (Gal et al., 2017); (2) diversity sampling methods such as Core-Set (Sener & Savarese, 2017) and Variational Adversarial Active Learning (VAAL) (Sinha et al., 2019); and (3) hybrid method such as Batch Active learning by Diverse Gradient Embeddings (BADGE) (Ash et al., 2019).

We test the above methods on the CIFAR10 (Krizhevsky et al., 2009) dataset starting with an initial labeled set with size $|\mathcal{L}_{init}| = 10$, and conducted 50 acquisition rounds where after each round B = 10 new samples are selected and labeled. We use ResNet-18 (He et al., 2016) as our training model across all methods. And we repeated the acquisitions five times with different random seeds. The results are visualized in Fig. 6, where we plot the *relative* performance between each method and random sampling acquisition through a diverging color map.



Figure 6: Existing methods fail to outperform random sampling with small budgets. This figure shows the relative performance between multiple methods and random acquisition. Within each subplot, x axis represents the accumulative acquisition size, while y axis indicates runs initiated with different random seeds. White color indicates on-par performance with random, blue indicates worse, and red indicates better.

800 Aligning with the general perceptions that low-budget (Mittal et al., 2019; Hacohen et al., 2022) and 801 cold-start (Zhu et al., 2019; Chandra et al., 2021) AL tasks are especially challenging, we empirically 802 observe that almost all popular AL algorithms fail to outperform the naive random sampling when 803 acquisition quota is less than 1% (500 out of 50,000 in the case of CIFAR10) of the unlabeled size. 804 More specifically, when the quota is less than 0.2% (less than 100 data points for CIFAR-10), all 805 methods fail to reliably outperform random sampling (as the beginning of each heatmap in Fig. 6 806 are almost all blue), which greatly motivates the development of DAO. We also include the more 807 conventional line plot of the empirical analysis which may provide more detailed information of each run in Fig. 7. 808



Figure 7: Relative performance between existing popular AL methods and random acquisition. horizontal axis represents the accumulative size of the labeled set, while vertical axis indicates relative performance in percentage.

B EXPERIMENTS ON DIFFERENT STOCHASTIC BATCH ACQUISITIONS

In this section, we provide more detailed results of §5.4. Specifically, we further study the performance of DAO when no batch sampling strategy or other sampling strategy is used and compare the results with existing popular AL algorithms. The results are shown in Fig. 8. For all experiments, we used the same low-budget setting as discussed in §4.



Figure 8: CIFAR-10 experiment results on (a): DAO without batch acquisition strategy (using naive top-k selection) and with other sampling strategies (softmax and softrank, as discussed in §3.5); (b): DAO without sampling (top-k) vs. existing AL algorithms; (c): DAO with softrank sampling vs. existing AL algorithms; (d): DAO with softmax sampling vs. existing AL algorithms; In all subplots, horizontal axis represents the accumulative size of the labeled set, while vertical axis indicates classification accuracy.

⁸⁶⁴ C VISUALIZATIONS OF SAMPLES SELECTED BY DAO

In this section, we show the visual representations of the data samples collected by DAO as a complement to §5.5. Unlabeled and newly acquired data, in this case, images, or their latent space embeddings, are first dimensionally-reduced and then visualized in Fig. 9. We see that, DAO-selected data exhibit characteristics of diversity across the sample space over multiple acquisition rounds, while display uncertainty characteristics within single round.



Figure 9: Visualizations of DAO acquisitions with dimensionality reduced from (a): raw images; and (b): latent space image embeddings.