LABELED TRUSTSET GUIDED: BATCH ACTIVE LEARNING WITH REINFORCEMENT LEARNING

Anonymous authors

000

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027 028 029

030

033

034

036

038

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Batch active learning (BAL) is a crucial technique for reducing labeling costs and improving data efficiency in training large-scale deep learning models. Traditional BAL methods often rely on metrics like Mahalanobis Distance to balance uncertainty and diversity when selecting data for annotation. However, these methods predominantly focus on the distribution of unlabeled data and fail to leverage feedback from labeled data or the model's performance. To address these limitations, we introduce TrustSet, a novel approach that selects the most informative data from the labeled dataset, ensuring a balanced class distribution to mitigate the long-tail problem. Unlike CoreSet, which focuses on maintaining the overall data distribution, TrustSet optimizes the model's performance by pruning redundant data and using label information to refine the selection process. To extend the benefits of TrustSet to the unlabeled pool, we propose a reinforcement learning (RL)-based sampling policy that approximates the selection of high-quality Trust-Set candidates from the unlabeled data. Combining TrustSet and RL, we introduce the Batch Reinforcement Active Learning with TrustSet (BRAL-T) framework. BRAL-T achieves state-of-the-art results across 10 image classification benchmarks and 2 active fine-tuning tasks, demonstrating its effectiveness and efficiency in various domains.

1 Introduction

In the era of deep learning, large-scale labeled datasets are indispensable for training models on complex tasks. Active learning (AL) provides an efficient approach to reduce the labeling costs by intelligently selecting critical subsets from unlabeled data for annotation (Zhan et al., 2022; Yang et al., 2024; Németh & Matuszka, 2024; Safaei et al., 2024). Batch active learning (BAL) (Citovsky et al., 2021), a variant of AL, further improves this process by selecting data points in groups (batches), thereby reducing the overhead associated with model retraining and oracle interactions.

In most modern BAL methods, the selection strategy is typically based on two factors: uncertainty and diversity. Uncertainty-based methods focus on choosing the most ambiguous or difficult data, which is likely to improve the model, but this often results in selecting redundant data that doesn't sufficiently cover the data distribution (Shen et al., 2017). On the other hand, diversity-based methods aim to ensure a representative subset by covering as many different types of data as possible, but they may neglect critical uncertain samples near the decision boundaries. For instance, CoreSet (Phillips, 2017) selects subsets that reflect the overall data distribution, ensuring diversity by minimizing the distance between the selected subset and the full dataset. While methods like Cluster-Margin (Citovsky et al., 2021) combine diversity and uncertainty to improve data selection, they still have limitations, such as overlooking feedback from the labeled dataset, ignoring class distribution, and potentially inheriting the long-tail distribution problem.

To address these challenges, we propose TrustSet, a novel data selection approach that distinguishes itself from CoreSet by emphasizing the utilization of label information. TrustSet focuses not only on ensuring diversity but also on selecting data that is most beneficial for improving the model's performance and releasing class imbalance problem. TrustSet differs from CoreSet in two ways:

Objective: TrustSet is designed to optimize the model's performance, with an explicit focus on improving accuracy and tackling the long-tail distribution problem by selecting crucial data that has a high potential to be forgotten by the model (Toneva et al., 2018). In contrast, CoreSet focuses

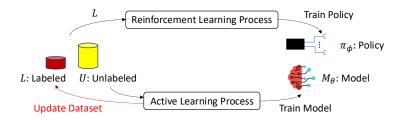


Figure 1: An overview of the BRAL-T framework, which consists of two main processes: (1)**Active Learning** for data selection. (2)**Reinforcement Learning** to approximate TrustSet selection.

on representing the full data distribution without directly considering the impact on the model's learning process, which can lead to inheriting undesirable distributional imbalances.

Data Source: TrustSet leverages labeled data, utilizing ground truth labels to prune redundant and noisy data and ensure that the selected subset is balanced across classes. This approach contrasts with CoreSet, which selects data purely from unlabeled pool, without considering feedback from trained model. As a result, CoreSet-based methods can miss the opportunity to incorporate critical information about model's current performance, potentially leading to suboptimal data selections.

TrustSet's balanced class distribution ensures better handling of the long-tail distribution problem, where underrepresented classes are more likely to be included in the training process. To construct TrustSet, we use the GradNd method (Paul et al., 2021), which ranks data based on the gradient norms of model updates, prioritizing data points that contribute most to model learning. Furthermore, to improve the data quality, we incorporate SuperLoss (Castells et al., 2020), which follows a curriculum learning strategy to assign higher importance to easier data in early training stages, while still considering difficult samples later.

However, extending TrustSet to the unlabeled data pool presents a challenge, as the selection process requires label information. To overcome this, we introduce an RL-based policy for approximating the selection of high-potential TrustSet candidates from the unlabeled data pool. Unlike previous RL-based active learning approaches, which often require frequent retraining of the model and rely heavily on complex reward structures (Fang et al., 2017; Zhang et al., 2023), our method minimizes retraining costs by leveraging TrustSet to guide the reinforcement learning process.

To this end, we propose a novel batch active learning framework called BRAL-T (Batch Reinforcement Active Learning with TrustSet extraction), which integrates TrustSet and RL-based policies for efficient data selection. The framework consists of two primary components: (1) Trust-Set extraction, which ensures that the labeled dataset contributes optimally to model performance and maintains a balanced class distribution, and (2) RL-based subset selection, where a learned policy selects from the unlabeled data pool to approximate TrustSet. This significantly reduces the need for repeated oracle queries and model retraining. As shown in Figure 1, BRAL-T is implemented with two processes: reinforcement learning (RL) for policy training and active learning (AL) for model training.

Our contributions are summarized as following: (1) We introduce TrustSet, a novel method for data selection that leverages label information to balance uncertainty, diversity, and class distribution, thus releasing class imbalanced issue and improving performance of AL. (2) We develop an RL-based data selection policy that bridges the gap between TrustSet's label dependency and the unlabeled setting of active learning, allowing for more efficient and targeted data selection. (3) We propose BRAL-T, a new batch active learning framework that integrates TrustSet and RL to reduce the computational burden of active learning while improving model performance. We demonstrate that BRAL-T achieves state-of-the-art performance across multiple image classification and active fine-tuning tasks.

2 Related Work

Active Learning: Active learning aims to reduce labeling costs by selecting informative data based on uncertainty and diversity. Shen et al. (2017) explored uncertainty-based methods for Named

Entity Recognition, later integrating diversity. Galaxy (Zhang et al., 2022) leveraged a model confidence graph to estimate uncertainty, while Yuan et al. (2020) applied self-supervised models for diversity-based selection. Recent works (Liu et al., 2019; Ash et al., 2019; Sinha et al., 2019; Margatina et al., 2021; Ash et al., 2021; Citovsky et al., 2021; Kim et al., 2021; Gentile et al., 2022) highlight the uncertainty-diversity trade-off, where uncertainty-focused methods risk redundancy, and diversity-based methods may overlook critical uncertain samples. However, the batch sampling methods in these works tend to ignore the distribution of selected data, leading to further redundancy

Batch Active Learning: Active learning methods, aiming to minimize oracle queries, prefer batch data processing over individual sample handling. Batch active learning approaches (Zhang et al., 2023; Ash et al., 2021; Kirsch et al., 2019; Citovsky et al., 2021; Sener & Savarese, 2017) concentrate on batch sampling to reduce costs and preserve subset distribution. BatchBald (Kirsch et al., 2019) addressed the lack of joint informativeness in batch sampling with an entropy-based selection. Cluster-Margin (Citovsky et al., 2021) used Hierarchical Agglomerative Clustering to identify highly uncertain samples at scale. Despite their success in computer vision tasks, existing methods overlook feedback from selected samples, such as accuracy changes, which could be crucial for refining data sampling strategies.

Active Learning with RL: RL has been explored to learn data selection policies (Zhang et al., 2023; Fang et al., 2017; Liu et al., 2019; Gong et al., 2022; Smit et al., 2021; Casanova et al., 2020). Some approaches(Fang et al., 2017; Gong et al., 2022; Smit et al., 2021; Casanova et al., 2020) defined rewards based on target model metrics (e.g., accuracy, AUROC), requiring frequent retraining and suffering from credit assignment issues. Instead, Liu et al. (2019) used the Mahalanobis distance for reward definition, while TAILOR (Zhang et al., 2023) formulated class balance as a reward signal. However, the former is task-specific (Re-ID), and the latter does not directly align with model accuracy. To address these limitations, we propose an RL-based active learning approach that avoids retraining while maintaining high correlation with task performance.

3 Problem Definition

In this section, we formally define the active learning problem in batch setting, following Sener & Savarese (2017), and introduce the TrustSet selection problem. We consider an \mathbb{C} -class classification task over a compact input space \mathcal{X} and a label space $\mathcal{Y}=\{1,\ldots,\mathbb{C}\}$. Our goal is to train a target model M_{θ} with parameters θ to optimize a loss function $\ell(M_{\theta}(\mathbf{x}),y):\mathcal{X}\times\mathcal{Y}\to R$, where $\hat{y}=M_{\theta}(\mathbf{x})$ refers to the predicted category by the target model and cross entropy loss is widely used as ℓ in classification task. In practice, we assume a large collection of data points sampled i.i.d. over the space $\mathcal{Z}=\mathcal{X}\times\mathcal{Y}$ as $\{\mathbf{x}_k,y_k\}_{k\in[n]}\sim p_{\mathcal{Z}}$, where n refers to the total amount of data points and $[n]=\{1,2,\ldots,n\}$.

For active learning problem, we further define the labeled dataset with |L| data points as $L=\{\mathbf{x}_k,y_k\}_k^{|L|}$ and the unlabeled data pool with |U| data points as $U=\{\mathbf{x}_k\}_k^{|U|}$. In general, $|L|\ll |U|$ and $L\cap U=\emptyset$. For the i-th active learning iteration, the labeled and unlabeled dataset are defined as L_i and U_i respectively. We aim to select a data subset $S_{U_i}\subseteq U_i$ to be labeled by an oracle and added to the current labeled dataset L_i forming an enhanced labeled dataset $L_{i+1}=L_i\cup S_{U_i}$. We abbreviate S_{U_i} as S_i for readability in the rest of the paper. Training on L_{i+1} , the model $M_{\theta_{L_{i+1}}}$ with the trained parameter $\theta_{L_{i+1}}$ is expected to achieve the best performance across all possible choices of S_i which can be formulated as the following Eq. 1:

$$S_i^* = \underset{S_i \subseteq U_i: |S_i| \le b}{\operatorname{arg \, min}} E_{\mathbf{x}, y \sim p_{\mathcal{Z}}} [\ell(M_{\theta_{L_{i+1}}}(\mathbf{x}), y)]$$
 (1)

where S_i^* refers to the optimal data subset for active learning and b refers to the size of data subset that needs to be selected in each iteration. As a result, Eq. 1 indicates that we aim to select a data subset S_i from U_i to enhance L_i such that the trained model $M_{\theta_{L_{i+1}}}$ achieves the minimal loss value.

Directly solving the above optimization problem is challenging due to the large number of possible choice of S_i from U_i . To address this, we consider an ideal scenario where label information for U_i is available and analyze the entire dataset $D = L_i \cup U_i$ to identify the most important samples that contribute to model training. Formally, we define the TrustSet T_D to be the important data of D as

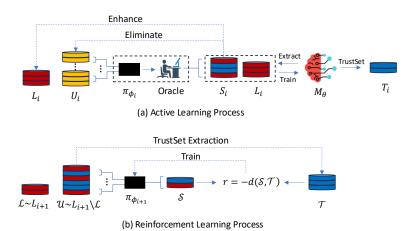


Figure 2: Details of BRAL-T. (a) Active Learning Process: π_{ϕ_i} selects the subset S_i from U_i for the oracle to annotate, and the model $M_{\theta_{L_{i+1}}}$ is trained on $S_i \cup L_i$. (b) RL process: we sample $\mathcal L$ and $\mathcal U$ from L_{i+1} to train policy $\pi_{\phi_{i+1}}$.

the following Eq. 2:

$$T_D = \underset{S_D \subseteq D: |S_D| \le b_T}{\arg \min} E_{\mathbf{x}, y \sim p_{\mathcal{Z}}} [\ell(M_{\theta_{S_D}}(\mathbf{x}), y)] \quad \text{s.t. balance}(S_D)$$
 (2)

where S_D refers to a data subset selected from D and we require S_D to be balanced across $\mathbb C$ categories for TrustSet selection to alleviate long-tail distribution problem; $M_{\theta_{S_D}}$ refers to the model trained on S_D and b_T refers to the predefined size of TrustSet. Eq. 2 indicates that given the limited size of data subset, T_D contains the most useful information from the dataset to train the model. Combining with Eq. 1, we instead select \tilde{S}_i from the unlabeled data pool U_i for active learning as:

$$\tilde{S}_i = \underset{S_i \subseteq U_i, |S_i| \le b}{\arg \min} d(S_i, T_D \cap U_i)$$
(3)

where $d(\cdot, \cdot)$ indicates a statistical distance function (e.g. Wasserstein Distance) for two distributions and we aim to select a data subset that has a distribution similar to $T_D \cap U_i$. As T_D contains the most significant data points for model training, enhancing L_i with such data points could also benefit the performance of the trained model. As a result, \tilde{S}_i , which approximates $T_D \cap U_i$, also contributes to a significant improvement of model training. Compared to Eq. 1, TrustSet is more reliable due to the use of annotation labels and the balance requirement releases the class imbalance problem of collected labeled dataset. Additionally, Eq. 2 can be more straightforwardly solved using data pruning methods (Paul et al., 2021; Park et al., 2023; Tan et al., 2023).

However, label information of U_i is not available under the active learning assumption. To solve this problem, we train a data selection policy π_{ϕ_i} with RL method on the labeled dataset L_i , learning to select \tilde{S}_i as an approximation of $T_D \cap U_i$. In the i-th iteration, we create an environment similar to the active learning setting by randomly sampling two subsets \mathcal{L} and \mathcal{U} from labeled dataset L_i , where the labels of \mathcal{L} are retained to simulate labeled dataset while the labels of \mathcal{U} are omitted to simulate unlabeled dataset. The whole dataset in this environment is then defined as $\mathcal{D} = \mathcal{L} \cup \mathcal{U}$, ensuring $\mathcal{L} \cap \mathcal{U} = \emptyset$ and $|\mathcal{L}| \ll |\mathcal{U}|$. Although we omit the labels of \mathcal{U} for active learning purpose, we can still leverage the label to extract $T_{\mathcal{D}}$. Taking \mathcal{L} and \mathcal{U} as input, π_{ϕ_i} is trained to select a subset \mathcal{S} from \mathcal{U} , with the reward defined as:

$$R = -d(\mathcal{S}, T_{\mathcal{D}} \cap \mathcal{U}) \tag{4}$$

where $S = \pi_{\phi_i}(\mathcal{L}, \mathcal{U})$ and we optimize the parameters ϕ_i to minimize the statistical distance between $T_D \cap \mathcal{U}$ and S. In this way, π_{ϕ_i} learns to select data from the unlabeled dataset that has a high potential to be included in the TrustSet based on the feature space. After training, π_{ϕ_i} is applied to the real unlabeled data pool U_i to select $S_i = \pi_{\phi_i}(L_i, U_i)$.

In general, for each active learning iteration, we solve Eq. 1 in two processes as shown in Figure 2. In the active learning process, data subset S_i is selected by $\pi_{\phi_i}(L_i, U_i)$ and passed to the oracle

for annotation. The model M_{θ} is then trained on the enhanced labeled dataset L_{i+1} . Moreover, the new unlabeled data pool U_{i+1} is achieved by eliminating S_i from U_i and the TrustSet T_i is extracted from L_{i+1} by solving Eq. 2 to facilitate the next RL process. In the RL process, we create the environment based on L_i and train π_{ϕ_i} using the reward function defined in Eq. 4.

4 METHOD

In this section, we introduce BRAL-T framework in detail. In Section 4.1, we introduce a TrustSet construction method based on GradNd score (Paul et al., 2021). In Section 4.2, we illustrate the details of the RL module and describe how we use the learned policy to select subsets from the unlabeled pool.

4.1 TRUSTSET

In general, the TrustSet should retain important data and tend to be class-balanced. However, it is almost impossible to directly solve Eq. 2 due to the large size of labeled dataset and time-consuming of $M_{\theta_{S_D}}$ training. As a result, we introduce a TrustSet extraction method based on the GradNd score (Paul et al., 2021) by analyzing the performance of model trained on entire trainset L rather than selected subset S. This score is defined as the expected value of the gradient norm term with respect to a differentiable model and a data sample x:

$$GradNd = E \| \sum_{k=1}^{K} \nabla_{M_{\theta}^{(k)}} \ell(M_{\theta}(x), y)^{T} \nabla_{\theta} M_{\theta}^{(k)}(x) \|$$
 (5)

In this equation, K denotes the number of logits, $M_{\theta}(x) \in \mathbb{R}^K$ refers to the output of model and $M_{\theta}^{(k)}(x)$ represents the result of the k-th logit from the model M_{θ} . For instance, in an image classification task, ℓ represents the cross-entropy loss, $K = \mathbb{C}$ is the number of categories, and $M_{\theta}^{(k)}(x)$ is the logit output for the k-th category. Data samples that result in a large gradient value tend to contain information that the model has not yet learned, as the model would update significantly based on such data. As demonstrated by the experimental analysis from Paul et al. (2021), data with a higher GradNd score tend to be forgotten samples for the target model during training and are more important for further training. However, the GradNd score might lead to a class imbalance problem when the data subset primarily contains difficult images for certain categories. To mitigate the long-tail distribution problem, we sort data by class using the GradNd score and select the top-N data for each category. For the image classification task, we follow Paul et al. (2021) in omitting the term $\nabla \theta M_{\theta}^{(k)}(x)$ from Eq 5 and calculate the EL2N score to approximate GradNd.

Curriculum Learning: Data with high GradNd scores tend to be difficult and uncertain samples. As suggested by previous works (Ash et al., 2021; Citovsky et al., 2021; Gentile et al., 2022), a training set focusing on uncertainty could result in high redundancy and fail to train a model that captures general features. We reconsider this issue from another important perspective. Difficult samples contain noise that can interfere with model predictions and increase the difficulty for the model to learn the boundaries between categories. With a limited amount of data, easy examples could help the model capture features and cluster data within the same category. Following the principles of curriculum learning (Tang & Huang, 2019; Castells et al., 2020), we assign larger weights to easier data samples in the early active learning iterations and leverage Super Loss (Castells et al., 2020) on top of the task loss ℓ . For each data sample (x, y), the super loss ℓ_s is defined as:

$$\ell_s(M_\theta(x), y) = (\ell(M_\theta(x), y) - \tau)\sigma + \lambda(\log \sigma)^2$$
(6)

where τ is the threshold for separating easy and hard samples, and λ is the weight of the regularization term. Both τ and λ are hyperparameters, while σ is learnable and indicates the weight assigned to the task loss. To minimize ℓ_s , data with task loss $\ell < \tau$ will be assigned a larger weight σ , and data with $\ell > \tau$ will be assigned a smaller weight σ . Since model training on uncertain data typically results in larger losses compared to easier data, super loss adaptively adjusts the weight for data samples. Meanwhile, the scale of σ is determined by λ . As λ increases, the value of σ tends to be 1 and has less effect on the task loss. Specifically, when $\lambda \to \infty$, σ will always be 1 to minimize the regularization term, making ℓ_s equal to $\ell - \tau$. As shown in Section 5.4, with Super Loss, proposed method achieves better performance. In the following sections, T_D refers to the TrustSet with Super Loss unless explicitly stated.

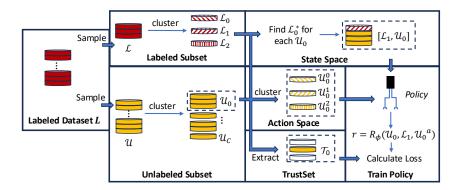


Figure 3: Reinforcement learning process. We use \mathcal{U}_0 as an example for process in the figure.

4.2 Reinforcement Learning

 The TrustSet is collected with label information to ensure class balance and improved reliability. However, during the active learning process, the policy needs to be applied to select a subset from the unlabeled data pool U_i . To address this, we create an environment with similar conditions and apply RL to train a policy for subset selection, where the TrustSet serves as the target subset. In the remainder of this section, we first define the state, action, and reward for the RL task in general, followed by an illustration of the overall process.

State $(\mathcal{L}_c^*, \mathcal{U}_c)$. We randomly sample \mathcal{L} as a labeled dataset and \mathcal{U} as an unlabeled data pool from L_i to train the policy π_{ϕ_i} . For each sampled dataset $\mathcal{D} = \mathcal{L} \cup \mathcal{U}$, we extract $T_{\mathcal{D}}$ and calculate $T_{\mathcal{D}} \cap \mathcal{U}$ as the target selected subset. For convenience, we define $\mathcal{T} = T_{\mathcal{D}} \cap \mathcal{U}$. Using all data in \mathcal{L} and \mathcal{U} as input is computationally expensive and challenging for learning an effective policy. It is beneficial to have alternative representations. Since in classification task, data tend to cluster based on predicted categories in feature space and \mathcal{T} is distributed across all clusters, it is more reasonable to predict \mathcal{T} using clusters as states. Thus, we define the state space as $(\mathcal{L}_c^*, \mathcal{U}_c)$, where \mathcal{U}_c refers to the c-th cluster from the unlabeled data pool \mathcal{U} , and \mathcal{L}_c^* is a cluster from the labeled dataset \mathcal{L} defined as:

$$\mathcal{L}_c^* = \operatorname*{arg\,min}_m d(\mathcal{L}_m, \mathcal{U}_c) \tag{7}$$

In this equation, \mathcal{L}_m refers to the m-th cluster from \mathcal{L} , and \mathcal{L}_c^* is the closest labeled cluster to \mathcal{U}_c based on the distance function $d(\cdot,\cdot)$. We use the Wasserstein Distance (Flamary et al., 2021) in our RL process. For each $(\mathcal{L},\mathcal{U})$ sample, there are C states, where C is the number of clusters in \mathcal{U} . For improved efficiency, we extract stochastic features of clusters as input for the policy, specifically using mean and variance as $[E[\mathcal{L}_c^*], Var[\mathcal{L}_c^*], E[\mathcal{U}_c], Var[\mathcal{U}_c]]$.

Action \mathcal{U}_c^a . As data in the TrustSet tend to group together by cluster, we further divide \mathcal{U}_c into A_c data groups, denoted as $\{\mathcal{U}_c^a\}_{a=1}^{A_c}$. Given $(\mathcal{L}_c^*,\mathcal{U}_c)$ as input, the policy selects the top clusters within $\{\mathcal{U}_c^a\}_{a=1}^{A_c}$ that have high potential to be included in the TrustSet. Consequently, $\{\mathcal{U}_c^a\}_{a=1}^{A_c}$ represents the candidate action space for each state, and the union of the selected actions forms the final selected data subset \mathcal{S} .

Reward R. Since different \mathcal{U}^a_c contain varying numbers of data points, for a fixed size of \mathcal{S} , we need to select a varying number of \mathcal{U}^a_c . It is more general to define the reward based on \mathcal{U}^a_c rather than \mathcal{S} . We set the reward as the negative Wasserstein distance between \mathcal{U}^a_c and the sub-TrustSet $\mathcal{T}_c = \mathcal{T} \cap \mathcal{U}_c$ as:

$$R = -d(\mathcal{U}_c^a, \mathcal{T}_c) \tag{8}$$

where \mathcal{U}_c^a closer in distribution to \mathcal{T}_c receives better reward.

We follow the DQN method (Mnih et al., 2013) and illustrate the overall RL process in Figure 3. Given the dataset \mathcal{L} and \mathcal{U} , we pass them through the target model M_{θ_i} to obtain the feature space. To construct the input state for the policy, we cluster the features of \mathcal{L} into M clusters, denoted as $\{\mathcal{L}_m\}_{m=1}^M$ and the features of \mathcal{U} into C clusters as $\{\mathcal{U}_c\}_{c=1}^C$. To generate candidate actions, we further divide \mathcal{U}_c into A_c clusters, denoted as $\{\mathcal{U}_c^a\}_{a=1}^{A_c}$. Meanwhile, we extract the TrustSet for each cluster as $\mathcal{T}_c = \mathcal{T} \cap \mathcal{U}_c$.

Unlike traditional RL tasks, we consider the future effect in curriculum learning and TrustSet selection, focusing only on the next timestep for policy training. As a result, training the Q-function is equivalent to training the reward function R_{ϕ} as:

$$\tilde{R} = R_{\phi}(\mathcal{U}_c, \mathcal{L}_c^*, \mathcal{U}_c^a) \tag{9}$$

where \tilde{R} refers to the predicted reward value by R_{ϕ} and ϕ refers to the parameter of the reward function. And ϕ is updated by the Mean Square Error (MSE) loss as:

$$l_{RL} = E_{(\mathcal{U}_c, \mathcal{U}_a^a)} \| R - \tilde{R} \| \tag{10}$$

During RL, we calculate the reward for all candidate actions and states to optimize the reward function R_{ϕ} . During the active learning process, for each unlabeled cluster U_c in the real unlabeled data pool U, we predict the reward for all candidate actions U_c^a and select them in descending order based on the reward score until the fixed size of subset selection is satisfied.

It is worthwhile to note that for each $(\mathcal{L}, \mathcal{U})$ sampled from L, M_{θ} needs to be retrained on $\mathcal{D} = \mathcal{L} \cup \mathcal{U}$ to extract the TrustSet $T_{\mathcal{D}}$ as well as \mathcal{T} , since the most important data for the model may vary depending on the labeled training set. To avoid the time-consuming process of frequent retraining, we approximate the extraction of $T_{\mathcal{D}}$ by reusing $M_{\theta_{L_i}}$ which has been trained on L_i in the i-th active learning iteration. This is justified by the fact that \mathcal{D} is randomly sampled from L_i and our main objective is to enhance L_i based on the performance of $M_{\theta_{L_i}}$ in the i-th active learning iteration. The only requirement is to retrain the policy from scratch for each active learning iteration. In practice, we extract T_i from the entire labeled set L_i and extract \mathcal{T} by taking the intersection between \mathcal{U} and T_i as $\mathcal{T} = T_i \cap \mathcal{U}$ for efficiency. Please refer to the Appendix for more details of RL training.

5 EXPERIMENT

In this section, we evaluate our proposed BRAL-T method on the image classification task and compare our results with previous active learning baselines, following the experimental settings of Zhan et al. (2022). Additionally, we also evaluate BRAL-T on the active fine-tuning task (Xie et al., 2023) and compare it with the current state-of-the-art method, ActiveFT. More experiments will be presented in Appendix.

5.1 IMAGE CLASSIFICATION RESULTS

Datasets: We evaluated BRAL-T on the image classification task across 8 benchmarks, including Cifar10, Cifar100 (Krizhevsky et al., 2009), Cifar10-imb, EMNIST (Cohen et al., 2017), Fashion-MNIST (Xiao et al., 2017), BreakHis (Spanhol et al., 2015), Pneumonia-MNIST (Kermany et al., 2018) and Waterbird (Sagawa et al., 2019; Koh et al., 2021). To create the Cifar10-imb dataset, we subsampled the training set of Cifar10 with ratios of 1:2:...:10 for classes 0 through 9.

Baselines: We compared BRAL-T with three baselines, LossPrediction (Yoo & Kweon, 2019), WAAL (Shui et al., 2020) and RandomSample. LossPrediction employs an additional module that predicts the loss for each data point. WAAL adopts min-max loss to better distinguish labeled and unlabeled samples while searching unlabeled batch with higher diversity than labeled samples. According to the experiments in Zhan et al. (2022), among all the methods, LossPrediction and WAAL achieve best results in 6 benchmarks and competitive results in other 2 benchmarks, therefore we select them as our baselines. For RandomSample, we randomly selected a subset from the unlabeled dataset in each active learning iteration. Besides, we visualized accuracy-budget curve on Cifar10, Cifar10-imb, Cifar100 and FashionMNIST benchmarks and compared with LossPrediction, WAAL, VAAL (Sinha et al., 2019), BADGE (Ash et al., 2019), CoreSet (Zhan et al., 2022), Cluster-Margin (Citovsky et al., 2021), BALD (Gal et al., 2017) and KMeans (Ash et al., 2019). To ensure a fair comparison, we used ResNet18 (He et al., 2016) as the target model. For more experimental details and hyperparameter settings, please refer to Appendix.

Evaluation Metrics: For all benchmarks, we report evaluation results using two metrics: *area under the budget curve* (AUBC) (Zhan et al., 2021a;b) and *final accuracy* (F-acc). AUBC refers to the area under the accuracy-budget curve. Methods with a higher AUBC score achieve better overall performance across different sizes of the training set. F-acc refers to the final accuracy achieved

Methods	Fashion	MNIST	EMN	IST	CIFA	R10	CIFA	R100
Methous	AUBC	F-acc	AUBC	F-acc	AUBC	F-acc	AUBC	F-acc
LossPrediction	0.859	0.888	0.762	0.793	0.837	0.911	0.481	0.655
WAAL	0.861	0.891	0.808	0.831	0.842	0.883	0.460	0.594
RandomSample	0.844	0.874	0.804	0.828	0.832	0.902	0.517	0.650
BRAL-T	0.863	0.894	0.813	0.833	0.847	0.916	0.525	0.662
	Cifar1	0-imb	Brea	kHis	Pneum.	MNIST	Wate	rbird
Benchmarks	AUBC	F-acc	AUBC	F-acc	AUBC	F-acc	AUBC	F-acc
LossPrediction	0.748	0.848	0.834	0.844	0.732	0.870	0.588	0.586
WAAL	0.752	0.799	0.836	0.855	0.640	0.870	0.525	0.506
WAAL RandomSample	0.752 0.710	0.799 0.810	0.836 0.834	0.855 0.832	0.640 0.706	0.870 0.652	0.525 0.586	0.506 0.502

Table 1: Experiment results of image classification task on 8 benchmarks.

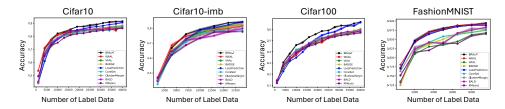


Figure 4: Visualization of experiment results on Cifar10, Cifar10-imb, Cifar100 and FashionMNIST.

after the budget Q is exhausted. The experiments for BRAL-T and the baselines were repeated for 3 trials under different random seeds, and the average of the evaluation results are reported.

Experiment Results: The experimental results on 8 benchmarks are presented in Table 1. BRALT significantly outperforms RandomSample under all benchmarks. Compared with WAAL and LossPrediction, BRAL-T achieves better AUBC as well as F-acc on all benchmarks. In Figure 4, we visualize the accuracy-budget curves of BRAL-T and baselines on 4 benchmarks. BRAL-T consistently achieves higher accuracy throughout the entire active learning process for Cifar100 and FashionMNIST. In early active learning iterations of Cifar10 and Cifar10-imb, BRAL-T has a bit worse accuracy compared with WAAL, the reasons of which could be attributed to WAAL's emphasis on diversity. However, without adequate consideration for uncertainty cause performance diminishing of WAAL when the size of labeled dataset increases. As comparison, LossPrediction focuses solely on uncertain data with high predicted loss, neglecting the diversity of the selected subset, which results in bad performance in early stage.

5.2 ACTIVE LEARNING ON MORE LONG-TAIL DATASETS

Dataset. Besides the aforementioned benchmarks, in this section, we focus on long-tail datasets, including CIFAR10-LT and CIFAR100-LT. Both datasets are subsampled from CIFAR datasets and the number of samples within each classes decreases exponentially with factor within 10 and 100. Specifically, we consider 10 and 20 in our experiments. The test images of CIFAR10-LT and CIFAR100-LT are the same as those in CIFAR10 and CIFAR100 datasets respectively. Both the two benchmarks are open-source and can be accessed through huggingface *tomas-gajarsky/cifar10-lt* and *tomas-gajarsky/cifar10-lt*. Please refer to Appendix for more details.

Experiment Results. Besides LossPrediction, WAAL and RandomSample, we also compare BRAL-T with SIMILAR (Kothawade et al., 2021) and TiDAL (Kye et al., 2023) which are designed for active learning in imbalanced dataset. As shown in Table 2, BRAL-T achieves the best AUBC as well as F-acc results. As pseudo label is not reliable especially when target model might be overconfident in long-tail dataset, BRAL-T selects informative data with ground-truth label to construct TrustSet which is more reliable to reflect whether target model has sufficiently learnt from related samples. As a result, BRAL-T always performs better than SIMILAR. Moreover, compared with LossPrediction and WAAL, we encourage TrustSet to be balanced which releases the category bias problem in long-tail distribution and contributes to the success of BRAL-T.

Methods	Cifar10	-LT-r10	Cifar10	-LT-r20	Cifar100)-LT-r10	Cifar100)-LT-r20
Methods	AUBC	F-acc	AUBC	F-acc	AUBC	F-acc	AUBC	F-acc
TiDAL	0.510	0.619	0.435	0.583	0.398	0.450	0.365	0.420
SIMILAR	0.507	0.654	0.431	0.566	0.405	0.514	0.366	0.460
LossPrediction	0.478	0.679	0.413	0.585	0.424	0.516	0.380	0.460
WAAL	0.510	0.623	0.435	0.589	0.381	0.458	0.342	0.412
RandomSample	0.476	0.686	0.397	0.587	0.382	0.463	0.343	0.410
BRAL-T	0.512	0.686	0.440	0.614	0.432	0.517	0.384	0.462

Table 2: Experiment results on Cifar10-LT and Cifar100-LT datasets.

76.0.1	Cifar	10-imb	TinyIr	nageNet
Methods	2%	3%	2%	$\overline{3}\%$
RandomSample ActiveFT BRAL-T	0.841	0.856	0.213	0.348
ActiveFT	0.838	0.851	0.289	0.359
BRAL-T	0.852	0.865	0.300	0.392

Baseline	Cifa		Cifar100		
Daseille	AUBC	F-acc	AUBC	F-acc	
PseudoScore	0.842	0.908	0.486	0.661	
BRAL-DiffSet	0.843	0.909	0.521	0.652	
BRAL-T w/o CL	0.845	0.906	0.522	0.662	
RandomSample	0.832	0.902	0.517	0.650	
BRAL-T	0.847	0.916	0.525	0.662	

Table 3: Results of Active Finetuning task.

Table 4: Ablation study result.

5.3 ACTIVE LEARNING FOR FINETUNING RESULTS

Experiment Setting: We adhere to the settings of Xie et al. (2023) and use Deit-Small (Touvron et al., 2021), pretrained with the DINO (Caron et al., 2021) framework on ImageNet-1k, as the target model. We chose two datasets for fine-tuning: Cifar10-imb and TinyImageNet (Le & Yang, 2015), resizing all images to 224×224 . For more implementation details, we utilize ActiveFT (Xie et al., 2023) to select an 1% subset as the initial labeled dataset and select an additional 1% of data for each active learning iteration. The pretrained model is fine-tuned using the SGD optimizer for 1000 epochs with a batch size of 512. Cosine learning rate decay is applied during the fine-tuning phase of each active learning iteration.

Experiment Results: All experiments were repeated across 3 trials, and the average results are reported in Table 3. BRAL-T significantly outperforms the other baselines. ActiveFT could suffer from the long-tail distribution of unlabeled data pool while TrustSet are defined to be class-balanced.

5.4 ABLATION STUDY

Experiment Setting: To further evaluate BRAL-T, we conducted ablation studies to demonstrate the benefits of the proposed modules by considering three baselines. For **PseudoScore**, instead of training an RL policy, we assign pseudo-labels to the unlabeled data pool based on the category with the highest logit score. Data subset with top EL2N score will be selected during active learning. For **BRAL-DiffSet**, to show the effectiveness of TrustSet, We select the second-best data group instead of the best as TrustSet. For **BRAL-T w/o CL**, we remove curriculum learning and directly use the cross-entropy loss function to calculate the EL2N score.

Experiment Results: Table 4 displays the results for the Cifar10 and Cifar100 datasets. BRAL-T surpasses PseudoScore in both AUBC and F-acc as pseudo-labels are often inaccurate, especially when classifier has low performance. In contrast, selecting the TrustSet based on the labeled dataset is more reliable. Compared to BRAL-DiffSet, BRAL-T also achieves better AUBC and F-acc scores, empirically proving the correlation between EL2N score and model accuracy. Moreover, curriculum learning also plays a crucial role in the success of BRAL-T, which aids in selecting easy examples and enhances the performance of the target model in the initial stages.

6 Conclusion

In summary, our RL-based Active Learning framework, BRAL-T, leverages TrustSet to more accurately evaluate distribution of labeled datasets and employs an RL policy to learn from Trust-Set. BRAL-T benchmarked against 8 baselines across 8 image classification tasks, shows superior AUBC and F-acc performance. Moreover, in CIFAR-LT benchmarks, BRAL-T outperforms baselines to handle long-tail dataset. Additionally, its application in active fine-tuning tasks reveals new state-of-the-art results.

REFERENCES

- Jordan Ash, Surbhi Goel, Akshay Krishnamurthy, and Sham Kakade. Gone fishing: Neural active learning with fisher embeddings. *Advances in Neural Information Processing Systems*, 34:8927–8939, 2021.
- 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.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
- Arantxa Casanova, Pedro O Pinheiro, Negar Rostamzadeh, and Christopher J Pal. Reinforced active learning for image segmentation. *arXiv preprint arXiv:2002.06583*, 2020.
- Thibault Castells, Philippe Weinzaepfel, and Jerome Revaud. Superloss: A generic loss for robust curriculum learning. *Advances in Neural Information Processing Systems*, 33:4308–4319, 2020.
- Gui Citovsky, Giulia DeSalvo, Claudio Gentile, Lazaros Karydas, Anand Rajagopalan, Afshin Rostamizadeh, and Sanjiv Kumar. Batch active learning at scale. *Advances in Neural Information Processing Systems*, 34:11933–11944, 2021.
- Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist to handwritten letters. In 2017 international joint conference on neural networks (IJCNN), pp. 2921–2926. IEEE, 2017.
- Meng Fang, Yuan Li, and Trevor Cohn. Learning how to active learn: A deep reinforcement learning approach. *arXiv preprint arXiv:1708.02383*, 2017.
- Rémi Flamary, Nicolas Courty, Alexandre Gramfort, Mokhtar Z. Alaya, Aurélie Boisbunon, Stanislas Chambon, Laetitia Chapel, Adrien Corenflos, Kilian Fatras, Nemo Fournier, Léo Gautheron, Nathalie T.H. Gayraud, Hicham Janati, Alain Rakotomamonjy, Ievgen Redko, Antoine Rolet, Antony Schutz, Vivien Seguy, Danica J. Sutherland, Romain Tavenard, Alexander Tong, and Titouan Vayer. Pot: Python optimal transport. *Journal of Machine Learning Research*, 22(78): 1–8, 2021. URL http://jmlr.org/papers/v22/20-451.html.
- 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.
- Claudio Gentile, Zhilei Wang, and Tong Zhang. Achieving minimax rates in pool-based batch active learning. In *International Conference on Machine Learning*, pp. 7339–7367. PMLR, 2022.
- Jia Gong, Zhipeng Fan, Qiuhong Ke, Hossein Rahmani, and Jun Liu. Meta agent teaming active learning for pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11079–11089, 2022.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pp. 630–645. Springer, 2016.
- Daniel S Kermany, Michael Goldbaum, Wenjia Cai, Carolina CS Valentim, Huiying Liang, Sally L Baxter, Alex McKeown, Ge Yang, Xiaokang Wu, Fangbing Yan, et al. Identifying medical diagnoses and treatable diseases by image-based deep learning. *cell*, 172(5):1122–1131, 2018.
- Kwanyoung Kim, Dongwon Park, Kwang In Kim, and Se Young Chun. Task-aware variational adversarial active learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8166–8175, 2021.
- Yoon-Yeong Kim, Youngjae Cho, JoonHo Jang, Byeonghu Na, Yeongmin Kim, Kyungwoo Song, Wanmo Kang, and Il-Chul Moon. Saal: sharpness-aware active learning. In *International Conference on Machine Learning*, pp. 16424–16440. PMLR, 2023.

- Andreas Kirsch, Joost Van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning. *Advances in neural information processing systems*, 32, 2019.
 - Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Irena Gao, et al. Wilds: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning*, pp. 5637–5664. PMLR, 2021.
 - Suraj Kothawade, Nathan Beck, Krishnateja Killamsetty, and Rishabh Iyer. Similar: Submodular information measures based active learning in realistic scenarios. *Advances in Neural Information Processing Systems*, 34:18685–18697, 2021.
 - Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
 - Seong Min Kye, Kwanghee Choi, Hyeongmin Byun, and Buru Chang. Tidal: Learning training dynamics for active learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 22335–22345, 2023.
 - Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. CS 231N, 7(7):3, 2015.
 - Zimo Liu, Jingya Wang, Shaogang Gong, Huchuan Lu, and Dacheng Tao. Deep reinforcement active learning for human-in-the-loop person re-identification. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 6122–6131, 2019.
 - Katerina Margatina, Giorgos Vernikos, Loïc Barrault, and Nikolaos Aletras. Active learning by acquiring contrastive examples. *arXiv preprint arXiv:2109.03764*, 2021.
 - Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
 - Gábor Németh and Tamás Matuszka. Compute-efficient active learning. *arXiv preprint arXiv:2401.07639*, 2024.
 - Dongmin Park, Seola Choi, Doyoung Kim, Hwanjun Song, and Jae-Gil Lee. Robust data pruning under label noise via maximizing re-labeling accuracy. *Advances in Neural Information Processing Systems*, 36:74501–74514, 2023.
 - Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet: Finding important examples early in training. *Advances in Neural Information Processing Systems*, 34:20596–20607, 2021.
 - Jeff M Phillips. Coresets and sketches. In *Handbook of discrete and computational geometry*, pp. 1269–1288. Chapman and Hall/CRC, 2017.
 - Bardia Safaei, VS Vibashan, Celso M de Melo, and Vishal M Patel. Entropic open-set active learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 4686–4694, 2024.
 - Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv* preprint arXiv:1911.08731, 2019.
 - Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. *arXiv preprint arXiv:1708.00489*, 2017.
 - Yanyao Shen, Hyokun Yun, Zachary C Lipton, Yakov Kronrod, and Animashree Anandkumar. Deep active learning for named entity recognition. *arXiv preprint arXiv:1707.05928*, 2017.
 - Changjian Shui, Fan Zhou, Christian Gagné, and Boyu Wang. Deep active learning: Unified and principled method for query and training. In *International Conference on Artificial Intelligence and Statistics*, pp. 1308–1318. PMLR, 2020.

- Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell. Variational adversarial active learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 5972–5981, 2019.
- Akshay Smit, Damir Vrabac, Yujie He, Andrew Y Ng, Andrew L Beam, and Pranav Rajpurkar. Medselect: Selective labeling for medical image classification combining meta-learning with deep reinforcement learning. *arXiv* preprint arXiv:2103.14339, 2021.
- Fabio A Spanhol, Luiz S Oliveira, Caroline Petitjean, and Laurent Heutte. A dataset for breast cancer histopathological image classification. *IEEE transactions on biomedical engineering*, 63 (7):1455–1462, 2015.
- Haoru Tan, Sitong Wu, Fei Du, Yukang Chen, Zhibin Wang, Fan Wang, and Xiaojuan Qi. Data pruning via moving-one-sample-out. *Advances in neural information processing systems*, 36: 18251–18262, 2023.
- Ying-Peng Tang and Sheng-Jun Huang. Self-paced active learning: Query the right thing at the right time. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pp. 5117–5124, 2019.
- Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J Gordon. An empirical study of example forgetting during deep neural network learning. *arXiv preprint arXiv:1812.05159*, 2018.
- Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In International conference on machine learning, pp. 10347–10357. PMLR, 2021.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- Yichen Xie, Han Lu, Junchi Yan, Xiaokang Yang, Masayoshi Tomizuka, and Wei Zhan. Active finetuning: Exploiting annotation budget in the pretraining-finetuning paradigm. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23715–23724, 2023.
- Chenhongyi Yang, Lichao Huang, and Elliot J Crowley. Plug and play active learning for object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 17784–17793, 2024.
- Donggeun Yoo and In So Kweon. Learning loss for active learning. In *Proceedings of the IEEE/CVF* conference on computer vision and pattern recognition, pp. 93–102, 2019.
- Michelle Yuan, Hsuan-Tien Lin, and Jordan Boyd-Graber. Cold-start active learning through self-supervised language modeling. *arXiv preprint arXiv:2010.09535*, 2020.
- Xueying Zhan, Qing Li, and Antoni B Chan. Multiple-criteria based active learning with fixed-size determinantal point processes. *arXiv preprint arXiv:2107.01622*, 2021a.
- Xueying Zhan, Huan Liu, Qing Li, and Antoni B Chan. A comparative survey: Benchmarking for pool-based active learning. In *IJCAI*, pp. 4679–4686, 2021b.
- Xueying Zhan, Qingzhong Wang, Kuan-hao Huang, Haoyi Xiong, Dejing Dou, and Antoni B Chan. A comparative survey of deep active learning. *arXiv preprint arXiv:2203.13450*, 2022.
- Jifan Zhang, Julian Katz-Samuels, and Robert Nowak. Galaxy: Graph-based active learning at the extreme. In *International Conference on Machine Learning*, pp. 26223–26238. PMLR, 2022.
- Jifan Zhang, Shuai Shao, Saurabh Verma, and Robert Nowak. Algorithm selection for deep active learning with imbalanced datasets. *arXiv preprint arXiv:2302.07317*, 2023.

A THE USAGE OF LARGE LANGUAGE MODEL (LLM)

LLM is only leveraged to polish the writing and correct wrong expressions in the paper. We provide each paragraph of the paper draft to ChatGPT separately and ask for revision. The core ideas, contributions of the paper and technique details are originated by the author without the LLM engagement.

B METHOD DETAILS

648

649 650

651

652

653

654 655

656 657

658

659

660

661

662

664

665

666

667 668

690

691

692

693

694

696

697

699

700

701

BRAL-T comprise two iterative processes: active learning process and reinforcement learning process. Algorithm 1 shows the pseudocode of overall framework. We randomly sampled initial labeled dataset and initialize parameters of target model and reward network in lines 3-5. During the i-th active learning process (lines 7-9), we trained target model M_{θ_i} with i-th labeled dataset L_i from scratch and extract TrustSet T_i from L_i , details of which is depicted in Section 4.1. During the reinforcement learning (RL) process (lines 11-17), we followed DQN (Mnih et al., 2013) and initialized replay buffer $\mathbb B$ to be empty. For each RL iteration, we sample labeled set and unlabeled set from L_i and store state set $\{\mathcal L_c^*, \mathcal U_c, \mathcal U_c^a, T_c\}$ into replay buffer (detailed in algorithm 2). To train reward function R_{ϕ_i} , a data batch is sampled from $\mathbb B$ and parameters of R_{ϕ_i} is updated based on Eq 10 (detailed in algorithm 3). After the two processes, we sampled a new dataset for oracle to annotate and updated L_i and U_i .

Algorithm 1 BRAL-T

```
669
            1: Input: Dataset D
670
            2: Output: Target Model M_{\theta}
671
            3: Random sample L_0 from D and annotated by oracle;
672
            4: Set U_0 := D \setminus L_0;
673
            5: Initialize M_{\theta_0}, R_{\phi_0};
674
            6: for i = 0 to N - 1 do
675
                   // Active Learning Process
            7:
676
                   Train M_{\theta_i} with L_i from scratch;
            8:
677
            9:
                   Extract TrustSet T_i from L_i;
678
           10:
           11:
                   // Reinforcement Learning Process
679
                   Initialize Replay Buffer B;
           12:
680
           13:
                   for j = 0 to K do
681
           14:
                       Sample \mathcal{L} and \mathcal{U} from L_i;
           15:
                      Extract set \{\mathcal{L}_c^*, \mathcal{U}_c, \mathcal{U}_c^a, \mathcal{T}_c\} = \mathbb{E}(\mathcal{L}, \mathcal{U}, T_i) and store into \mathbb{B}; (Algorithm 2)
683
           16:
                       Sample data from \mathbb{B} and train R_{\phi_i} as Eq 10. (Algorithm 3)
684
           17:
                   end for
685
           18:
686
                   // Sample New DataSet
           19:
687
                   Sample S_i := \pi(R_{\phi_i}, L_i, U_i);
                   Update L_{i+1} := L_i \cup S_i and U_{i+1} := U_i \setminus S_i;
688
           21:
           22: end for
689
```

In algorithm 2, we show the pseudocode of data extraction for RL (line 15 of algorithm 1). As illustrated in Section 4.2, we clustered labeled set into $\{L_m\}_{m=1}^M$ and unlabeled set into $\{U_c\}_{c=1}^C$ to formulate state space of RL. For each unlabeled subset, we further cluster U_c into $\{U_c^a\}_{a=1}^{A_c}$ to formulate action space of RL and extract Trustset T_c for each U_c . All pairs of $\{L_c^*, U_c, U_c^a, T_c\}$ are stored and return as extraction results.

In algorithm 3, we show the pseudocode of reinforcement learning to train data selection policy (line 16 of algorithm 1). For each gradient step, we sample state and action data from replay buffer $\mathbb B$ and extract vector input S and A (line 4-6). Then based on Eq. 8, we calculate the reward for each (state, action) pair as negative distance between data subset and TrustSet (line 7). And based on Eq. 9, we predict reward with current reward function R_{ϕ} (line 8). Finally, we calculate mean square error (MSE) loss between predicted reward r and ground truth reward R and update reward function with gradient descent (line 9).

Algorithm 2 Data Extraction For Reinforcement Learning

```
703
                 1: Input: LabeledSet \mathcal{L}, UnlabeledSet \mathcal{U}, TrustSet T
704
                 2: Output: Data list Out
705
                 3: Initialize output list Out := [];
706
                 4: Cluster \mathcal{L} into \{\mathcal{L}_m\}_{m=1}^M;
5: Cluster \mathcal{U} into \{\mathcal{U}_c\}_{c=1}^C
708
                 6: for each U_c do
                           Extract \mathcal{T}_c := T_i \cap \mathcal{U}_c
Cluster \mathcal{U}_c into \{\mathcal{U}_c^a\}_{a=1}^{A_c};
709
                 7:
710
                 8:
                 9:
                           Calculate \mathcal{L}_c^* := \arg\min_m d(\mathcal{L}_m, \mathcal{U}_c);
711
                           Store each \{\mathcal{L}_{c}^{*}, \mathcal{U}_{c}, \mathcal{U}_{c}^{a}, \mathcal{T}_{c}\} into Out;
                10:
712
               11: end for
713
               12: Return Out;
714
```

Algorithm 3 Training of Reinforcement Learning

```
    Input: Replay Buffer B, Reward Function R<sub>φ</sub>.
    Output: Update Reward Function R<sub>φ</sub>
    for Each Gradient Step do
    Sample data batch from B as {L<sub>c</sub>*, U<sub>c</sub>, U<sub>c</sub>*, T<sub>c</sub>}<sup>B</sup>.
    Extract state vector input as: S = [E[L<sub>c</sub>*], Var[L<sub>c</sub>*], E[U<sub>c</sub>], Var[U<sub>c</sub>]].
    Extract action input as: A = [E[U<sub>c</sub>*], Var[U<sub>c</sub>*]].
    Calculate reward for each state action pair as Eq. 8: R = -d(A, T<sub>c</sub>).
    Predicate reward with R<sub>φ</sub> as Eq. 9: r = R<sub>φ</sub>(S, A).
    Calculate Loss L = MSE(R, r) and update R<sub>φ</sub> with gradient descent.
```

10: **end for** 11: Return R_{ϕ} ;

C EXPERIMENT DETAILS

In this section, we introduce more experiment details of Section 5, including architecture of target model we used for image classification and hyperparameter settings of experiments.

Benchmarks	$ L_0 $	$ U_0 $	Q	b	#e	\overline{C}
FashionMNIST	500	59,500	10,000	250	40	10
EMNIST	1,000	696,932	50,000	500	40	62
CIFAR10	1,000	49,000	40,000	500	50	10
CIFAR100	1,000	49,000	40,000	500	60	100
CIFAR10-imb	1,000	27,239	20,000	500	50	10
CIFAR10-LT	1,000	-	D	100	50	10
CIFAR100-LT	4,000	-	D	500	60	100
BreakHis	100	5,436	5,000	100	30	2
PneumoniaMNIST	100	5,132	5,000	100	30	2
Waterbird	100	4,695	4,000	100	30	2

Table 5: Setting of benchmarks. Where $|L_0|$ refers to size of initial labeled set, $|U_0|$ refers to size of initial unlabeled data pool, Q refers to budget, b refers to batch size for target model training, #e refers to number of epoch for target model training and C refers to number of clusters from unlabeled data pool. For all the benchmarks, the number of clusters M from labeled dataset is set to be the same as C and number of candidate action for U_c is set to be 5.

DataSets. We evaluated BRAL-T on the image classification task across 5 benchmarks, including Cifar10, Cifar100 (Krizhevsky et al., 2009), Cifar10-imb, EMNIST (Cohen et al., 2017), and FashionMNIST (Xiao et al., 2017). To create the Cifar10-imb dataset, we followed the settings of Zhan et al. (2022) and subsampled the training set with ratios of 1:2:...:10 for classes 0 through 9. We also evaluated our framework on medical imaging analysis tasks across 2 benchmarks, including Breast cancer Histopathological Image Classification (BreakHis) (Spanhol et al., 2015) and Chest X-Ray

Pneumonia classification (Pneumonia-MNIST) (Kermany et al., 2018). Additionally, we assessed our framework on an object recognition dataset with correlated backgrounds (Waterbird) (Sagawa et al., 2019; Koh et al., 2021), which contains waterbird and landbird classes manually mixed with water and land backgrounds. To further evaluate BRAL-T on long-tail datasets, we also consider CIFAR10-LT and CIFAR100-LT where the number of samples within each classes decreases exponentially with factor to be 10, 20 or 50.

The detail setting for each benchmark are shown in Table 5, including initial data size of labeled dataset $|L_0|$ and unlabeled dataset $|U_0|$, final budget Q of labeled dataset, batch size b for data subset selection in each active learning iteration, training epoch #e for target model training, and category number C for dataset.

Model Details. Following the setting of Zhan et al. (2022), we use Resnet18 (He et al., 2016) as the target model for image classification tasks. For Cifar10, Cifar10-imb, Cifar100 and PneumoniaMNIST, we replaced the kernal size of first convolutional layer to be 3×3 and stride to be 1 in order to handle image with smaller size. For grayscale images such as FashionMNIST and EMNIST datasets, we add an additional convolutional layer before the first layer of Resnet with 1×1 kernal to increase the channel number of images to be 3. Furthermore, we trained the target models of all baselines for the same number of epochs, as shown in Table 5. For LossPrediction, the target model is trained with both classification loss and loss prediction loss for the first 20 epochs. After 20 epochs, only the gradient from the classification loss is back-propagated through the target model.

Model and Hyperparameters Setting: We constructed the reward function R_{ϕ} using a fully connected network comprising 2 hidden layers, each with 512 units, and use the ReLU activation function. SGD was employed as the optimizer for R_{ϕ} , with the learning rate set at 0.01. For hyperparameters of curriculum learning, we follow the setting of SuperLoss ((Castells et al., 2020)) and set $\tau = \log |K|$ where |K| is the category number. Additionally, we set the value of λ to be 0.25 for EMNIST, CIFAR100 and TinyImageNet datasets and 1.0 for the others. During active learning, We train target model with SGD optimizer for PneumoniaMNIST and Waterbird benchmarks and Adam optimizer for other datasets.

After each active learning iteration, we sampled 30 pairs of \mathcal{L} and \mathcal{U} from the existing labeled set L to train the policy, setting the batch size to 100 pairs of state, action, and reward. Following each sampling, we trained R_{ϕ} for 20 iterations, resulting in a total of 600 iterations for the entire RL training process. As shown in Table 5, the number of clusters C for unlabeled set and M for labeled set are set to be the same as category number for related benchmark. And the number of candidate action A_c for each unlabeled cluster U_c is set to be 5 during the experiment.

D MORE EXPERIMENT RESULTS

In this section, we introduce more experiments and results. First of all, we show the confidence interval results for Table 1 over 8 benchmarks in D.1. Then we evaluate BRAL-T by calculating penalty matrix in D.3. Moreover, to show the efficiency of BRAL-T, we compare time overhead between BRAL-T and baselines in D.2. Finally, in D.4, we show more ablation studies of BRAL-T.

D.1 CONFIDENCE INTERVALS OF RESULTS IN IMAGE CLASSIFICATION TASKS.

Besides representing average value of AUBC and F-acc of BRAL-T and baselines on image classif-cation benchmarks, Table 6 shows the confidence interval of experiment results. In general, BRAL-T results are stable and robust over different experiment trials.

D.2 TIME OVERHEAD COMPARISON

To evaluate the efficiency of BRAL-T, we compare the time overhead with LossPrediction, WAAL, VAAL and SIMILAR on Cifar10 and Cifar100 datasets. All experiments were conducted using a single Quadro RTS 6000 GPU core with CUDA Version 11.4. Figure 5 shows the time cost results along with active learning iteration.

The time cost associated with BRAL-T increases with each active learning iteration as the labeled set expands and more samples are clustered during the reinforcement learning process. However, com-

Methods	Fashion	MNIST	EMN	NIST	CIFA	AR10	CIFA	R100
Methous	AUBC	F-acc	AUBC	F-acc	AUBC	F-acc	AUBC	F-acc
		± 0.038						
WAAL	± 0.002	± 0.015	± 0.012	± 0.015	± 0.006	± 0.009	± 0.006	± 0.011
RandomSample	± 0.001	± 0.009	± 0.004	± 0.007	± 0.003	± 0.011	± 0.003	± 0.008
BRAL-T	± 0.001	± 0.008	± 0.005	± 0.014	± 0.003	± 0.006	± 0.004	± 0.009
Dl	Cifar	10-imb	Brea	kHis	Pneum.	MNIST	Water	bird
Benchmarks	1	10-imb F-acc						
Benchmarks LossPrediction	AUBC		AUBC	F-acc	AUBC	F-acc	AUBC	F-acc
	AUBC ±0.011	F-acc	AUBC ±0.026	F-acc ±0.037	AUBC ±0.023	F-acc ±0.038	AUBC ±0.014	F-acc ±0.097
LossPrediction	#0.011 ±0.008	F-acc ±0.017 3 ±0.013	AUBC ±0.026 ±0.016	F-acc ±0.037 ±0.042	AUBC ±0.023 ±0.018	F-acc ±0.038 ±0.021	AUBC ±0.014 : ±0.011 :	F-acc ±0.097 ±0.078

Table 6: Confidence Interval of Experiment results of image classification task.

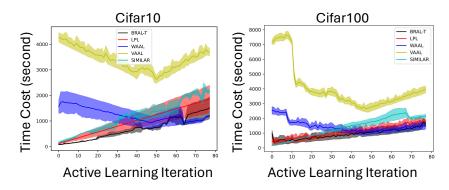


Figure 5: Time Cost of BRAL-T and baselines.

pared to other baselines, BRAL-T consistently demonstrates efficiency, maintaining a competitive edge in terms of computational resource utilization.

D.3 PAIRWISE COMPARISON

We further compare BRAL-T with VAAL (Sinha et al., 2019), SAAL (Kim et al., 2023) and BAIT (Ash et al., 2021) on Cifar10, Cifar10-imb and FashionMNIST datasets by pairwise penalty matrix following Ash et al. (2021). For each benchmark, we collect accuracy results achieved by all baselines. For pairwise comparison between the method for ith row (r_i) and the method in jth column (c_j), we add a score to element e_{ij} whenever r_i achieves better accuracy result in one budget of data subset for a benchmark, which means the better r_i performs compared with c_j , the higher score e_{ij} will be.

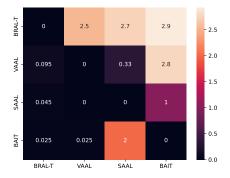


Figure 6: Pairwise Comparison of BRAL-T, VAAL, SAAL and BAIT.

Figure 6 represents the pairwise comparison results. Compared with all baselines, BRAL-T achieves highest value in e_{BRAL-T} , and lowest value in e_{BRAL-T} .

Agent Reuse. A potential way to further improve the efficiency of BRAL-T is reusing RL agent for all active learning iterations. However, considering the distribution shift of labeled dataset, distribution of TrustSet will also shift during active learning. For this reason, we apply two RL agents during active learning, one of which is trained in the first active learning step and remains unchanged for early active learning iterations; the other one of which is maintained for the rest iterations. Specifically for CIFAR10-imb dataset, we use the first agent for the first 20 iterations and the second agent for the rest 20 iterations. The result is shown in Table 7 below: where BRAL-T

Method	AUBC	F-Acc
LossPrediction	0.748	0.848
WAAL	0.752	0.799
RandomSample	0.710	0.810
BRAL-T	0.762	0.851
BRAL-T (two agents)	0.755	0.837

Table 7: BRAL-T reusing two RL agents.

with agent reusing surprisingly achieves better AUBC results compares with other baselines. With a more careful separation of active learning stages and RL agents, we believe the performance could be further improved.

D.4 MORE ABLATION STUDY

To evaluate the robustness of BRAL-T, we run BRAL-T on Cifar10-imb dataset under different qualifies of initial labeled dataset. Moreover, we show the performance of BRAL-T with different candidate action numbers. To evaluate the quality of RL approximation, we apply ground truth labels for TrustSet selection and compare the accuracy results with BRAL-T.

Quality Effect of Initial Labeled Set. We explored the impact of the initial labeled set's quality by applying three different sampling methods to construct the initial labeled set from the Cifar10-imb dataset:

- **Random Sample**: We randomly sample data from the unlabeled pool to form the initial labeled set which maintains a similar category distribution with unlabeled pool.
- Twisted Main: We sort the 10 categories by the number of data samples first and then select 50 samples from 5 rare classes and 950 samples randomly from the other 5 main classes.
- **Twisted Rare**: Similar to Twisted Main, we randomly select 50 samples from 5 main classes and 950 samples from the other 5 rare classes.

The results, depicted in the Figure 7, indicate that BRAL-T's performance varies with the quality of the initial labeled set, particularly when labeled data is scarce. However, as the size of the labeled dataset increases, the accuracy differences become negligible, demonstrating BRAL-T's robustness to the initial set's composition. Despite the initial set's quality impacting BRAL-T's performance, in Table 8, the AUBC results in the twisted cases are competitive with the results of the WAAL baseline in Table-2, and all achieve better F-Acc compared with other baselines.

Ablation Study on Different Action Numbers. In the reinforcement learning process, we set number of candidate action to be 5 in Section 5. To evaluate the impact of varying action space sizes,

Initial Method	AUBC	F-Acc
Random	0.762	0.851
Twisted Main	0.750	0.855
Twisted Rare	0.756	0.855

Table 8: Experiment Result of Different Initial Labeled Set Quality.

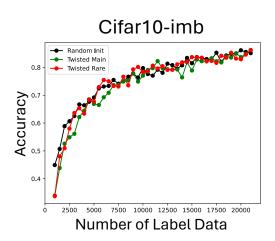


Figure 7: Accuracy-budget curve of Different Initial Labeled Set Quality.

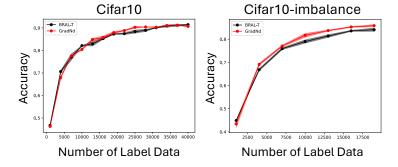


Figure 8: Comparison between BRAL-T with RL policy and Ground Truth Labels.

 we conducted an ablation study on the Cifar10-imb dataset, comparing BRAL-T's performance across different numbers of actions: 5, 10, 50, and 100. The results are shown in Table 9

# Actions	AUBC	F-Acc
5	0.762	0.851
10	0.763	0.853
50	0.755	0.851
100	0.758	0.854

Table 9: Ablation Study on Different Candidate Action Number.

Under all different setting of action numbers, BRAL-T achieves best AUBC and F-Acc results compared with baselines in Table 1. Setting a large number of actions will increase the complexity of policy training. As we keep the policy architecture to be the same and simple for time efficiency, in some active learning iteration policy might not be trained well with large action number which lead to a small drop of AUBC score. But in general, our method is robust to action number. The reason we choose 5 in the experiment is mainly for the consideration of time efficiency.

Ablation Study on Different Setting of λ **.** During the TrustSet extraction, we introduce curriculum learning where λ is introduced to control the effect of SuperLoss. We study the impact of λ on the CIFAR10-imb dataset for further sensitivity analysis, the result is shown in table 10.

Increasing the value of λ reduces the influence of SuperLoss on the task loss. In an imbalanced dataset, data samples are limited, especially in rare classes. Focusing on difficult data during the early stages of active learning can significantly increase the difficulty of model training. As a result, increasing λ leads to a reduction in AUBC and F-Acc for BRAL-T, highlighting the importance of incorporating curriculum learning into the active learning process. However, overall, the AUBC and F-Acc values remain competitive with the baselines presented in Table 1 of the paper.

Compare between RL and Ground Truth Labels. Although label information of unlabeled data pool is not available during active learning, in order to evaluate the approximation performance of RL policy, for baseline GradNd we assume ground truth label of unlabeled data pool is available when calculating the GradNd score of data samples and we pick class-balanced data with top GradNd score for each active learning iteration. We compare BRAL-T with GradNd on Cifar10 and Cifar10-imb datasets and shows the accuracy-budget results in Figure 8.

In Cifar10 dataset, BRAL-T achieves good performance to approximate TrustSet, where only small gap exists when labeled dataset becomes larger. In Cifar10-imb dataset, similarly, when labeled dataset is limited, BRAL-T achieves similar accuracy compared with GradNd. When the size of labeled dataset becomes larger, the accuracy difference performs to be acceptable larger. As a conclusion, the RL policy in BRAL-T achieves good performance to approximate ground truth TrustSet selection.

λ	AUBC	F-Acc
0.25	0.762	0.851
1.00	0.752	0.842
2.00	0.749	0.830

Table 10: Impact of λ value on CIFAR10-imb dataset.