# Making Your First Choice: To Address Cold Start Problem in Vision Active Learning

Liangyu Chen[1]    Yutong Bai[2]    Siyu Huang[3]    Yongyi Lu[2]
Bihan Wen[1]    Alan L. Yuille[2]    Zongwei Zhou[2],*
[1]Nanyang Technological University    [2]Johns Hopkins University    [3]Harvard University

## Abstract

Active learning promises to improve annotation efficiency by iteratively selecting the most important data to be annotated first. However, we uncover a striking contradiction to this promise: active learning fails to select data as efficiently as random selection at the first few choices. We identify this as the cold start problem in vision active learning, caused by a biased and outlier initial query. This paper seeks to address the cold start problem by exploiting the three advantages of contrastive learning: (1) no annotation is required; (2) label diversity is ensured by pseudo-labels to mitigate bias; (3) typical data is determined by contrastive features to reduce outliers. Experiments are conducted on CIFAR-10-LT and three medical imaging datasets (*i.e.* Colon Pathology, Abdominal CT, and Blood Cell Microscope). Our initial query not only significantly outperforms existing active querying strategies but also surpasses random selection by a large margin. We foresee our solution to the cold start problem as a simple yet strong baseline to choose the initial query for vision active learning.

Code is available: https://github.com/c-liangyu/CSVAL

## 1 Introduction

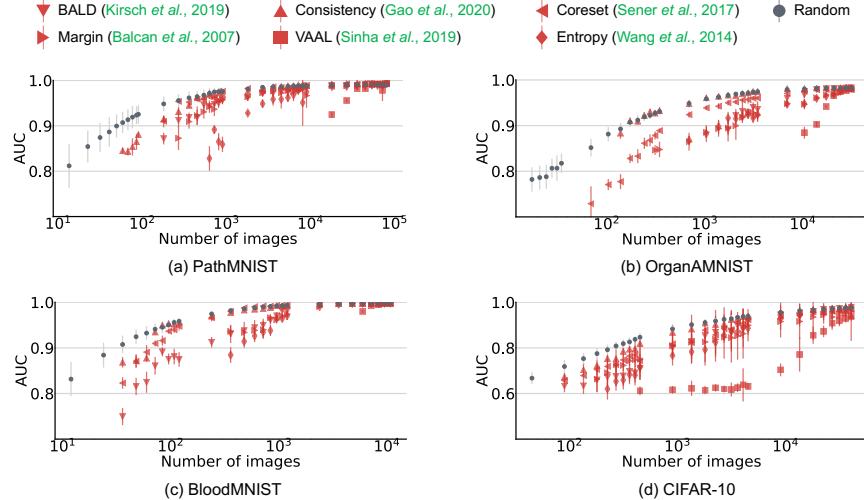*"The secret of getting ahead is getting started."*

— Mark Twain

The cold start problem was initially found in recommender systems [56, 39, 9, 23] when algorithms had not gathered sufficient information about users with no purchase history. It also occurred in many other fields, such as natural language processing [55, 33] and computer vision [5, 11, 38] during the active learning procedure[1]. Active learning promises to improve annotation efficiency by iteratively selecting the most important data to annotate. However, we uncover a striking contradiction to this promise: Active learning fails to select data as effectively as random selection at the first choice. We identify this as the cold start problem in vision active learning and illustrate the problem using three medical imaging applications (Figure 1a–c) as well as a natural imaging application (Figure 1d). Cold start is a crucial topic [54, 30] because a performant initial query can lead to noticeably improved subsequent cycle performance in the active learning procedure, evidenced in §3.3. There is a lack of studies that systematically illustrate the cold start problem, investigate its causes, and provide practical solutions to address it. To this end, we ask: *What causes the cold start problem and how can we select the initial query when there is no labeled data available?*

---

*Corresponding author: Zongwei Zhou (zzhou82@jh.edu)

[1]Active learning aims to select the most important data from the unlabeled dataset and query human experts to annotate new data. The newly annotated data is then added to improve the model. This process can be repeated until the model reaches a satisfactory performance level or the annotation budget is exhausted.

Figure 1: **Cold start problem in vision active learning.** Most existing active querying strategies (*e.g.* BALD, Consistency, etc.) are outperformed by random selection in selecting initial queries, since random selection is i.i.d. to the entire dataset. However, some classes are not selected by active querying strategies due to selection bias, so their results are not presented in the low budget regime.

Random selection is generally considered a baseline to start the active learning because the randomly sampled query is independent and identically distributed (i.i.d.) to the entire data distribution. As is known, maintaining a similar distribution between training and test data is beneficial, particularly when using limited training data [25]. Therefore, a large body of existing work selects the initial query randomly [10, 61, 55, 62, 18, 17, 42, 24, 22, 60], highlighting that active querying compromises accuracy and diversity compared to random sampling at the beginning of active learning [36, 63, 44, 11, 20, 59]. Why? We attribute the causes of the cold start problem to the following two aspects:

(i) *Biased query*: Active learning tends to select data that is biased to specific classes. Empirically, Figure 2 reveals that the class distribution in the selected query is highly unbalanced. These active querying strategies (*e.g.* Entropy, Margin, VAAL, etc.) can barely outperform random sampling at the beginning because some classes are simply not selected for training. It is because data of the minority classes occurs much less frequently than those of the majority classes. Moreover, datasets in practice are often highly unbalanced, particularly in medical images [32, 58]. This can escalate the biased sampling. We hypothesize that the *label diversity* of a query is an important criterion to determine the importance of the annotation. To evaluate this hypothesis theoretically, we explore the upper bound performance by enforcing a uniform distribution using ground truth (Table 1) To evaluate this hypothesis practically, we pursue the label diversity by exploiting the pseudo-labels generated by $K$-means clustering (Table 2). The label diversity can reduce the redundancy in the selection of majority classes, and increase the diversity by including data of minority classes.

(ii) *Outlier query*: Many active querying strategies were proposed to select typical data and eliminate outliers, but they heavily rely on a trained classifier to produce predictions or features. For example, to calculate the value of Entropy, a trained classifier is required to predict logits of the data. However, there is no such classifier at the start of active learning, at which point no labeled data is available for training. To express informative features for reliable predictions, we consider contrastive learning, which can be trained using unlabeled data only. Contrastive learning encourages models to discriminate between data augmented from the same image and data from different images [15, 13]. Such a learning process is called instance discrimination. We hypothesize that instance discrimination can act as an alternative to select typical data and eliminate outliers. Specifically, the data that is hard to discriminate from others could be considered as typical data. With the help of Dataset Maps [48, 26][2], we evaluate this hypothesis and propose a novel active querying strategy that can effectively select *typical data* (*hard-to-contrast* data in our definition, see §2.2) and reduce outliers.

---

[2]It is worthy noting that both [48] and [26] conducted a retrospective study, which analyzed existing active querying strategies by using the ground truth. As a result, the values of *confidence* and *variability* in the Dataset

Systematic ablation experiments and qualitative visualizations in §3 confirm that (i) the level of label diversity and (ii) the inclusion of typical data are two explicit criteria for determining the annotation importance. Naturally, contrastive learning is expected to approximate these two criteria: pseudo-labels in clustering implicitly enforce label diversity in the query; instance discrimination determines typical data. Extensive results show that our initial query not only significantly outperforms existing active querying strategies, but also surpasses random selection by a large margin on three medical imaging datasets (*i.e.* Colon Pathology, Abdominal CT, and Blood Cell Microscope) and two natural imaging datasets (*i.e.* CIFAR-10 and CIFAR-10-LT). Our active querying strategy eliminates the need for manual annotation to ensure the label diversity within initial queries, and more importantly, starts the active learning procedure with the typical data.

To the best of our knowledge, we are among the first to indicate and address the cold start problem in the field of medical image analysis (and perhaps, computer vision), making three contributions: (1) illustrating the cold start problem in vision active learning, (2) investigating the underlying causes with rigorous empirical analysis and visualization, and (3) determining effective initial queries for the active learning procedure. Our solution to the cold start problem can be used as a strong yet simple baseline to select the initial query for image classification and other vision tasks.

**Related work.** When the cold start problem was first observed in recommender systems, there were several solutions to remedy the insufficient information due to the lack of user history [63, 23]. In natural language processing (NLP), Yuan *et al.* [55] were among the first to address the cold start problem by pre-training models using self-supervision. They attributed the cold start problem to model instability and data scarcity. Vision active learning has shown higher performance than random selection [61, 47, 18, 2, 43, 34, 62], but there is limited study discussing how to select the initial query when facing the entire unlabeled dataset. A few studies somewhat indicated the existence of the cold start problem: Lang *et al.* [30] explored the effectiveness of the $K$-center algorithm [16] to select the initial queries. Similarly, Pourahmadi *et al.* [38] showed that a simple $K$-means clustering algorithm worked fairly well at the beginning of active learning, as it was capable of covering diverse classes and selecting a similar number of data per class. Most recently, a series of studies [20, 54, 46, 37] continued to propose new strategies for selecting the initial query from the entire unlabeled data and highlighted that typical data (defined in varying ways) could significantly improve the learning efficiency of active learning at a low budget. In addition to the existing publications, our study justifies the two causes of the cold start problem, systematically presents the existence of the problem in six dominant strategies, and produces a comprehensive guideline of initial query selection.

## 2 Method

In this section, we analyze in-depth the cause of cold start problem in two perspectives, biased query as the inter-class query and outlier query as the intra-class factor. We provide a complementary method to select the initial query based on both criteria. §2.1 illustrates that label diversity is a favourable selection criterion, and discusses how we obtain label diversity via simple contrastive learning and $K$-means algorithms. §2.2 describes an unsupervised method to sample atypical (hard-to-contrast) queries from Dataset Maps.

### 2.1 Inter-class Criterion: Enforcing Label Diversity to Mitigate Bias

$K$**-means clustering.** The selected query should cover data of diverse classes, and ideally, select similar number of data from each class. However, this requires the availability of ground truth, which are inaccessible according to the nature of active learning. Therefore, we exploit pseudo-labels generated by a simple $K$-means clustering algorithm and select an equal number of data from each cluster to form the initial query to facilitate label diversity. Without knowledge about the exact number of ground-truth classes, over-clustering is suggested in recent works [51, 57] to increase performances on the datasets with higher intra-class variance. Concretely, given 9, 11, 8 classes in the ground truth, we set $K$ (the number of clusters) to 30 in our experiments.

**Contrastive features.** $K$-means clustering requires features of each data point. Li *et al.* [31] suggested that for the purpose of clustering, contrastive methods (*e.g.* MoCo, SimCLR, BYOL) are

---

Maps could not be computed under the practical active learning setting because the ground truth is a priori unknown. Our modified strategy, however, does not require the availability of ground truth (detailed in §2.2).
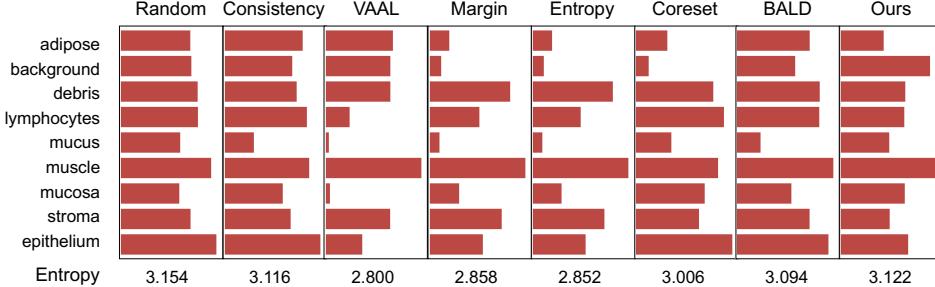
Figure 2: **Label diversity of querying criteria.** Random, the leftmost strategy, denotes the class distribution of randomly queried samples, which can also reflect the approximate class distribution of the entire dataset. As seen, even with a relatively larger initial query budget (40,498 images, 45% of the dataset), most active querying strategies are biased towards certain classes in the PathMNIST dataset. For example, VAAL prefers selecting data in the muscle class, but largely ignores data in the mucus and mucosa classes. On the contrary, our querying strategy selects more data from minority classes (e.g., mucus and mucosa) while retaining the class distribution of major classes. Similar observations in OrganAMNIST and BloodMNIST are shown in Appendix Figure 7. The higher the entropy is, the more balanced the class distribution is.

more suitable than generative methods (*e.g.* colorization, reconstruction) because the contrastive feature matrix can be naturally regarded as cluster representations. Therefore, we use MoCo v2 [15]— a popular self-supervised contrastive method—to extract image features.

$K$-means and MoCo v2 are certainly not the only choices for clustering and feature extraction. We employ these two well-received methods for simplicity and efficacy in addressing the cold start problem. Figure 2 shows our querying strategy can yield better label diversity than other six dominant active querying strategies; similar observations are made in OrganAMNIST and BloodMNIST (Figure 7) as well as CIFAR-10 and CIFAR-10-LT (Figure 10).

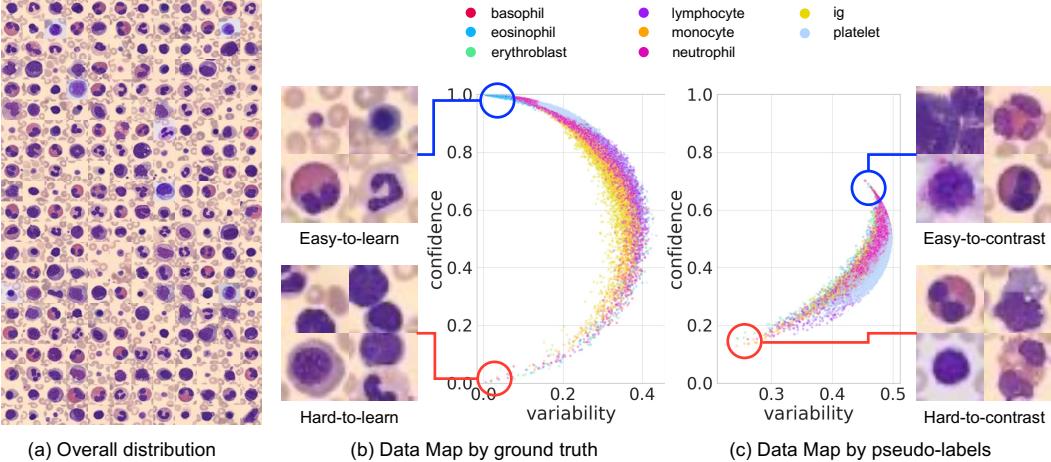## 2.2 Intra-class Criterion: Querying Hard-to-Contrast Data to Avoid Outliers

**Dataset map.** Given $K$ clusters generated from Criterion #1, we now determine which data points ought to be selected from each cluster. Intuitively, a data point can better represent a cluster distribution if it is harder to contrast itself with other data points in this cluster—we consider them typical data. To find these typical data, we modify the original Dataset Map[3] by replacing the ground truth term with a pseudo-label term. This modification is made because ground truths are unknown in the active learning setting but pseudo-labels are readily accessible from Criterion #1. For a visual comparison, Figure 3b and Figure 3c present the Data Maps based on ground truths and pseudo-labels, respectively. Formally, the modified Data Map can be formulated as follows. Let $\mathcal{D} = \{\boldsymbol{x}_m\}_{m=1}^{M}$ denote a dataset of $M$ unlabeled images. Considering a minibatch of $N$ images, for each image $\boldsymbol{x}_n$, its two augmented views form a positive pair, denoted as $\tilde{\boldsymbol{x}}_i$ and $\tilde{\boldsymbol{x}}_j$. The contrastive prediction task on pairs of augmented images derived from the minibatch generate $2N$ images, in which a true label $y_n^*$ for an anchor augmentation is associated with its counterpart of the positive pair. We treat the other $2(N-1)$ augmented images within a minibatch as negative pairs. We define the probability of positive pair in the instance discrimination task as:

$$p_{i,j} = \frac{\exp(\text{sim}(\boldsymbol{z}_i, \boldsymbol{z}_j))/\tau}{\sum_{n=1}^{2N} \mathbb{1}_{[n \neq i]} \exp(\text{sim}(\boldsymbol{z}_i, \boldsymbol{z}_n))/\tau}, \tag{1}$$

$$p_{\theta(e)}(y_n^*|x_n) = \frac{1}{2}[p_{2n-1,2n} + p_{2n,2n-1}], \tag{2}$$

where $\text{sim}(\boldsymbol{u}, ) = \boldsymbol{u}^\top \boldsymbol{v}/\|\boldsymbol{u}\|\|\boldsymbol{v}\|$ is the cosine similarity between $\boldsymbol{u}$ and $\boldsymbol{v}$; $\boldsymbol{z}_{2n-1}$ and $\boldsymbol{z}_{2n}$ denote the projection head output of a positive pair for the input $\boldsymbol{x}_n$ in a batch; $\mathbb{1}_{[n \neq i]} \in \{0, 1\}$ is an indicator

---

[3]Dataset Map [12, 48] was proposed to analyze datasets by two measures: *confidence* and *variability*, defined as the mean and standard deviation of the model probability of ground truth along the learning trajectory.

Figure 3: **Active querying based on Dataset Maps.** (a) Dataset overview. (b) Easy- and hard-to-learn data can be selected from the maps based on ground truths [26]. This querying strategy has two limitations: it requires manual annotations and the data are stratified by classes in the 2D space, leading to a poor label diversity in the selected queries. (c) Easy- and hard-to-contrast data can be selected from the maps based on pseudo-labels. This querying strategy is label-free and the selected hard-to-contrast data represent the most common patterns in the entire dataset, as presented in (a). These data are more suitable for training, and thus alleviate the cold start problem.

function evaluating to 1 iff $n \neq i$ and $\tau$ denotes a temperature parameter. $\theta^{(e)}$ denotes the parameters at the end of the $e^{\text{th}}$ epoch. We define confidence $(\hat{\mu}_m)$ across $E$ epochs as:

$$\hat{\mu}_m = \frac{1}{E} \sum_{e=1}^{E} p_{\theta^{(e)}}(y_m^*|x_m). \tag{3}$$

The confidence $(\hat{\mu}_m)$ is the Y-axis of the Dataset Maps (see Figure 3b-c).

**Hard-to-contrast data.** We consider the data with a low confidence value (Equation 3) as "hard-to-contrast" because they are seldom predicted correctly in the instance discrimination task. Apparently, if the model cannot distinguish a data point with others, this data point is expected to carry typical characteristics that are shared across the dataset [40]. Visually, hard-to-contrast data gather in the bottom region of the Dataset Maps and "easy-to-contrast" data gather in the top region. As expected, hard-to-learn data are more typical, possessing the most common visual patterns as the entire dataset; whereas easy-to-learn data appear like outliers [54, 26], which may not follow the majority data distribution (examples in Figure 3a and Figure 3c). Additionally, we also plot the original Dataset Map [12, 48] in Figure 3b, which grouped data into hard-to-learn and easy-to-learn[4]. Although the results in §3.2 show equally compelling performance achieved by both easy-to-learn [48] and hard-to-contrast data (ours), the latter do not require any manual annotation, and therefore are more practical and suitable for vision active learning.

In summary, to meet the both criteria, our proposed active querying strategy includes three steps: (i) extracting features by self-supervised contrastive learning, (ii) assigning clusters by $K$-means algorithm for label diversity, and (iii) selecting hard-to-contrast data from dataset maps.

## 3 Experimental Results

**Datasets & metrics.** Active querying strategies have a selection bias that is particularly harmful in long-tail distributions. Therefore, unlike most existing works [38, 54], which tested on highly balanced annotated datasets, we deliberately examine our method and other baselines on long-tail datasets to simulate real-world scenarios. Three medical datasets of different modalities

---

[4]Swayamdipta *et al.* [48] indicated that easy-to-learn data facilitated model training in the low budget regime because easier data reduced the confusion when the model approaching the rough decision boundary. In essence, the advantage of easy-to-learn data in active learning aligned with the motivation of curriculum learning [6].

Table 1: **Diversity is a significant add-on to most querying strategies.** AUC scores of different querying strategies are compared on three medical imaging datasets. In either low budget (*i.e.* 0.5% or 1% of MedMNIST datasets) or high budget (*i.e.* 10% or 20% of CIFAR-10-LT) regimes, both random and active querying strategies benefit from enforcing the label diversity of the selected data. The cells are highlighted in blue when adding diversity performs no worse than the original querying strategies. Coreset [41] works very well as its original form because this querying strategy has implicitly considered the label diversity (also verified in Table 2) by formulating a $K$-center problem, which selects $K$ data points to represent the entire dataset. Some results are missing (marked as "-") because the querying strategy fails to sample at least one data point for each class. Results of more sampling ratios are presented in Appendix Figures 6, 9.
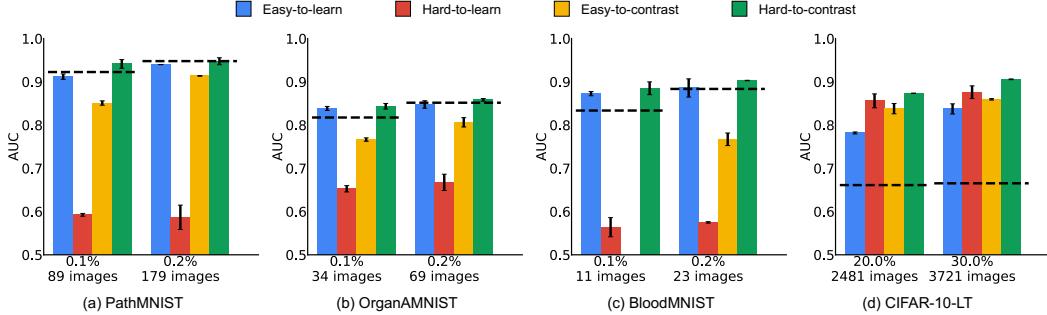
| | Unif. | PathMNIST 0.5% (499) | PathMNIST 1% (899) | OrganAMNIST 0.5% (172) | OrganAMNIST 1% (345) | BloodMNIST 0.5% (59) | BloodMNIST 1% (119) | CIFAR-10-LT 10% (1420) | CIFAR-10-LT 20% (2841) |
|---|---|---|---|---|---|---|---|---|---|
| Random | ✓ | 96.8±0.6 | 97.6±0.6 | 91.1±0.9 | 93.3±0.4 | 94.7±0.7 | 96.5±0.4 | 91.6±1.1 | 93.1±0.6 |
| | ✗ | 96.4±1.3 | 97.6±0.9 | 90.7±1.1 | 93.1±0.7 | 93.2±1.5 | 95.8±0.7 | 62.0±6.1 | - |
| Consistency | ✓ | 96.4±0.1 | 97.9±0.1 | 92.3±0.5 | 92.8±1.0 | 92.9±0.9 | 95.9±0.5 | 91.4±1.1 | 93.4±0.2 |
| | ✗ | 96.2±0.0 | 97.6±0.0 | 91.0±0.3 | 94.0±0.6 | 87.9±0.2 | 95.5±0.5 | 67.1±17.1 | 88.6±0.3 |
| VAAL | ✓ | 92.7±0.5 | 93.0±0.6 | 70.6±1.9 | 84.6±0.5 | 89.8±1.3 | 93.4±0.9 | 92.6±0.2 | 93.7±0.4 |
| | ✗ | - | - | - | - | - | - | - | - |
| Margin | ✓ | 97.9±0.2 | 96.0±0.4 | 81.8±1.2 | 85.8±1.4 | 89.7±1.9 | 94.7±0.7 | 91.7±0.9 | 93.2±0.2 |
| | ✗ | 91.0±2.3 | 96.0±0.3 | - | 85.9±0.7 | - | - | 81.9±0.8 | 86.3±0.3 |
| Entropy | ✓ | 93.2±1.6 | 95.2±0.2 | 79.1±2.3 | 86.7±0.8 | 85.9±0.5 | 91.8±1.0 | 92.0±1.2 | 91.9±1.3 |
| | ✗ | - | 87.5±0.1 | - | - | - | - | 65.6±15.6 | 86.4±0.2 |
| Coreset | ✓ | 95.0±2.2 | 94.8±2.5 | 85.6±0.4 | 89.9±0.5 | 88.5±0.6 | 94.1±1.1 | 91.5±0.4 | 93.6±0.2 |
| | ✗ | 95.6±0.7 | 97.5±0.2 | 83.8±0.6 | 88.5±0.4 | 87.3±1.6 | 94.0±1.2 | 65.9±15.9 | 86.9±0.1 |
| BALD | ✓ | 95.8±0.2 | 97.0±0.1 | 87.2±0.3 | 89.2±0.3 | 89.9±0.8 | 92.7±0.7 | 92.8±0.1 | 90.8±2.4 |
| | ✗ | 92.0±2.3 | 95.3±1.0 | - | - | 83.3±2.2 | 93.5±1.3 | 64.9±14.9 | 84.7±0.6 |

Table 2: **Class coverage of selected data.** Compared with random selection (i.i.d. to entire data distribution), most active querying strategies contain selection bias to specific classes, so the class coverage in their selections might be poor, particularly using low budgets. As seen, using 0.002% or even smaller proportion of MedMNIST datasets, the class coverage of active querying strategies is much lower than random selection. By integrating $K$-means clustering with contrastive features, our querying strategy is capable of covering 100% classes in most scenarios using low budgets ($\leq$0.002% of MedMNIST). We also found that our querying strategy covers the most of the classes in the CIFAR-10-LT dataset, which is designatedly more imbalanced.

| | PathMNIST 0.00015% (13) | PathMNIST 0.00030% (26) | OrganAMNIST 0.001% (34) | OrganAMNIST 0.002% (69) | BloodMNIST 0.001% (11) | BloodMNIST 0.002% (23) | CIFAR-10-LT 0.2% (24) | CIFAR-10-LT 0.3% (37) |
|---|---|---|---|---|---|---|---|---|
| Random | **0.79±0.11** | 0.95±0.07 | 0.91±0.08 | 0.98±0.04 | 0.70±0.13 | 0.94±0.08 | 0.58±0.10 | 0.66±0.12 |
| Consistency | 0.78 | 0.88 | 0.82 | 0.91 | 0.75 | 0.88 | 0.50 | 0.70 |
| VAAL | 0.11 | 0.11 | 0.18 | 0.18 | 0.13 | 0.13 | 0.30 | 0.30 |
| Margin | 0.67 | 0.78 | 0.73 | 0.82 | 0.63 | 0.75 | 0.60 | 0.70 |
| Entropy | 0.33 | 0.33 | 0.45 | 0.73 | 0.63 | 0.63 | 0.40 | 0.70 |
| Coreset | 0.66 | 0.78 | 0.91 | 1.00 | 0.63 | 0.88 | 0.60 | 0.70 |
| BALD | 0.33 | 0.44 | 0.64 | 0.64 | 0.75 | 0.88 | 0.60 | 0.70 |
| Ours | 0.78 | **1.00** | **1.00** | **1.00** | **1.00** | **1.00** | **0.70** | **0.80** |

in MedMNIST [53] are used: PathMNIST (colorectal cancer tissue histopathological images), BloodMNIST (microscopic peripheral blood cell images), OrganAMNIST (axial view abdominal CT images of multiple organs). OrganAMNIST is augmented following Azizi *et al*. [3], while the others following Chen *et al*. [15]. Area Under the ROC Curve (AUC) and Accuracy are used as the evaluation metrics. All results were based on at least three independent runs, and particularly, 100 independent runs for random selection. UMAP [35] is used to analyze feature clustering results.

**Baselines & implementations.** We benchmark a total of seven querying strategies: (1) random selection, (2) Max-Entropy [52], (3) Margin [4], (4) Consistency [18], (5) BALD [28], (6) VAAL [45], and (7) Coreset [41]. For contrastive learning, we trained 200 epochs with MoCo v2, following its default hyperparameter settings. We set $\tau$ to 0.05 in equation 2. To reproduce the large batch size and iteration numbers in [13], we apply repeated augmentation [21, 49, 50] (detailed in Table 5). More baseline and implementation details can be found in Appendix A.

Figure 4: **Quantitative comparison of map-based querying strategies.** Random selection (dot-lines) can be treated as a highly competitive baseline in cold start because it outperforms six popular active querying strategies as shown in Figure 1. In comparison with random selection and three other querying strategies, hard-to-contrast performs the best. Although easy-to-learn and hard-to-learn sometimes performs similarly to hard-to-contrast, their selection processes require ground truths [26], which are not available in the setting of active learning.

## 3.1 Contrastive Features Enable Label Diversity to Mitigate Bias

**Label coverage & diversity.** Most active querying strategies have selection bias towards specific classes, thus the class coverage in their selections might be poor (see Table 2), particularly at low budgets. By simply enforcing label diversity to these querying strategies can significantly improve the performance (see Table 1), which suggests that the label diversity is one of the causes that existing active querying strategies perform poorer than random selection.

Our proposed active querying strategy, however, is capable of covering 100% classes in most low budget scenarios ($\leq$0.002% of full dataset) by integraing $K$-means clustering with contrastive features.

## 3.2 Pseudo-labels Query Hard-to-Contrast Data and Avoid Outliers

**Hard-to-contrast data are practical for cold start problem.** Figure 4 presents the quantitative comparison of four map-based querying strategies, wherein easy- or hard-to-learn are selected by the maps based on ground truths, easy- or hard-to-contrast are selected by the maps based on pseudo-labels. Note that easy- or hard-to-learn are enforced with label diversity, due to their class-stratified distributions in the projected 2D space (illustrated in Figure 3). Results suggest that *selecting easy-to-learn or hard-to-contrast data contribute to the optimal models*. In any case, easy- or hard-to-learn data can not be selected without knowing ground truths, so these querying strategies are not practical for active learning procedure. Selecting hard-to-contrast, on the other hand, is a label-free strategy and yields the highest performance amongst existing active querying strategies (reviewed in Figure 1). More importantly, hard-to-contrast querying strategy significantly outperforms random selection by 1.8% (94.14%±1.0% vs. 92.27%±2.2%), 2.6% (84.35%±0.7% vs. 81.75%±2.1%), and 5.2% (88.51%±1.5% vs. 83.36%±3.5%) on PathMNIST, OrganAMNIST, and BloodMNIST, respectively, by querying 0.1% of entire dataset. Similarly on CIFAR-10-LT, hard-to-contrast significantly outperforms random selection by 21.2% (87.35%±0.0% vs. 66.12%±0.9%) and 24.1% (90.59%±0.1% vs. 66.53%±0.5%) by querying 20% and 30% of entire dataset respectively. Note that easy- or hard-to-learn are not enforced with label diversity, for a more informative comparison.

## 3.3 On the Importance of Selecting Superior Initial Query

**A good start foresees improved active learning.** We stress the importance of the cold start problem in vision active learning by conducting correlation analysis. Starting with 20 labeled images as the initial query, the training set is increased by 10 more images in each active learning cycle. Figure 14a presents the performance along the active learning (each point in the curve accounts for 5 independent trials). The initial query is selected by a total of 9 different strategies[5], and subsequent queries are

---

[5]Hard-to-learn is omitted because it falls behind other proposed methods by a large margin (Figure 4).
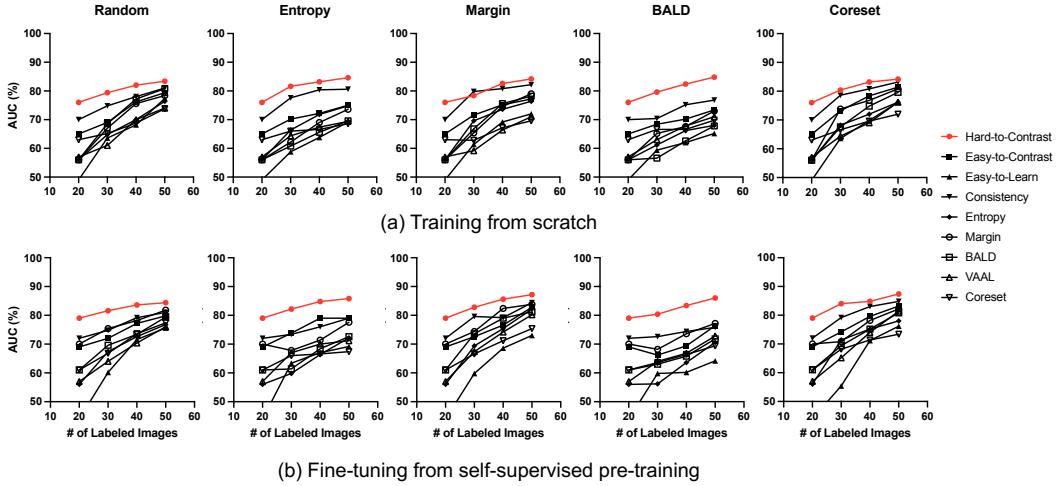
Figure 5: **On the importance of selecting a superior initial query.** Hard-to-contrast data (red lines) outperform other initial queries in every cycle of active learning on OrganaMNIST. We find that the performance of the initial cycle (20 images) and the last cycle (50 images) are strongly correlated.

selected by 5 different strategies. $AUC_n$ denotes the AUC score achieved by the model that is trained by $n$ labeled images. The Pearson correlation coefficient between $AUC_{20}$ (starting) and $AUC_{50}$ (ending) shows strong positive correlation ($r = 0.79, 0.80, 0.91, 0.67, 0.92$ for random selection, Entropy, Margin, BALD, and Coreset, respectively). This result is statistically significant ($p < 0.05$). Hard-to-contrast data (our proposal) consistently outperforms the others on OrganAMNIST (Figure 5), BloodMNIST (Figure 13), and PathMNIST (Figure 14), and steadily improves the model performances within the next active learning cycles.

**The initial query is consequential regardless of model initialization.** A pre-trained model can improve the performance of each active learning cycle for both random and active selection [55], but the cold start problem remains (evidenced in Figure 14b). This suggests that the model instability and data scarcity are two independent issues to be addressed for the cold start problem. Our "hard-to-contrast" data selection criterion only exploits contrastive learning (an improved model), but also determines the typical data to be annotated first (a better query). As a result, when fine-tuning from MoCo v2, the Pearson correlation coefficient between $AUC_{20}$ and $AUC_{50}$ remains high ($r = 0.92$, $0.81, 0.70, 0.82, 0.85$ for random selection, Entropy, Margin, BALD, and Coreset, respectively) and statistically significant ($p < 0.05$).

## 4 Conclusion

This paper systematically examines the causes of the cold start problem in vision active learning and offers a practical and effective solution to address this problem. Analytical results indicate that (1) the level of label diversity and (2) the inclusion of hard-to-contrast data are two explicit criteria to determine the annotation importance. To this end, we devise a novel active querying strategy that can enforce label diversity and determine hard-to-contrast data. The results of three medical imaging and two natural imaging datasets show that our initial query not only significantly outperforms existing active querying strategies but also surpasses random selection by a large margin. This finding is significant because it is the first few choices that define the efficacy and efficiency of the subsequent learning procedure. We foresee our solution to the cold start problem as a simple, yet strong, baseline to sample the initial query for active learning in image classification.

**Limitation.** This study provides an empirical benchmark of initial queries in active learning, while more theoretical analyses can be provided. Yehuda *et al.* [54] also found that the choice of active learning strategies depends on the initial query budget. A challenge is to articulate the quantity of determining active learning strategies, which we leave for future work.

**Potential societal impacts.** Real-world data often exhibit long-tailed distributions, rather than the ideal uniform distributions over each class. We improve active learning by enforcing label diversity and hard-to-contrast data. However, we only extensively test our strategies on academic datasets. In many other real-world domains such as robotics and autonomous driving, the data may impose additional constraints on annotation accessibility or learning dynamics, e.g., being fair or private. We focus on standard accuracy and AUC as our evaluation metrics while ignoring other ethical issues in imbalanced data, especially in underrepresented minority classes.

## Acknowledgements

## References

[1] Andrea Acevedo, Anna Merino, Santiago Alférez, Ángel Molina, Laura Boldú, and José Rodellar. A dataset of microscopic peripheral blood cell images for development of automatic recognition systems. *Data in Brief*, 30, 2020.

[2] Sharat Agarwal, Himanshu Arora, Saket Anand, and Chetan Arora. Contextual diversity for active learning. *ArXiv*, abs/2008.05723, 2020.

[3] Shekoofeh Azizi, Basil Mustafa, Fiona Ryan, Zachary Beaver, Jan Freyberg, Jonathan Deaton, Aaron Loh, Alan Karthikesalingam, Simon Kornblith, Ting Chen, et al. Big self-supervised models advance medical image classification. *arXiv preprint arXiv:2101.05224*, 2021.

[4] Maria-Florina Balcan, Andrei Broder, and Tong Zhang. Margin based active learning. In *International Conference on Computational Learning Theory*, pages 35–50. Springer, 2007.

[5] Javad Zolfaghari Bengar, Joost van de Weijer, Bartlomiej Twardowski, and Bogdan Raducanu. Reducing label effort: Self-supervised meets active learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1631–1639, 2021.

[6] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *ICML '09*, 2009.

[7] Maxim Berman, Hervé Jégou, Andrea Vedaldi, Iasonas Kokkinos, and Matthijs Douze. Multigrain: a unified image embedding for classes and instances. *ArXiv*, abs/1902.05509, 2019.

[8] Patrick Bilic, Patrick Ferdinand Christ, Eugene Vorontsov, Grzegorz Chlebus, Hao Chen, Qi Dou, Chi-Wing Fu, Xiao Han, Pheng-Ann Heng, Jürgen Hesser, et al. The liver tumor segmentation benchmark (lits). *arXiv preprint arXiv:1901.04056*, 2019.

[9] Jesús Bobadilla, Fernando Ortega, Antonio Hernando, and Jesús Bernal. A collaborative filtering approach to mitigate the new user cold start problem. *Knowledge-based systems*, 26:225–238, 2012.

[10] Alexander Borisov, Eugene Tuv, and George Runger. Active batch learning with stochastic query by forest. In *JMLR: Workshop and Conference Proceedings (2010)*. Citeseer, 2010.

[11] 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*, pages 14–32. PMLR, 2021.

[12] Haw-Shiuan Chang, Erik Learned-Miller, and Andrew McCallum. Active bias: Training more accurate neural networks by emphasizing high variance samples. *Advances in Neural Information Processing Systems*, 30, 2017.

[13] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020.

[14] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Moco demo: Cifar-10. `https://colab.research.google.com/github/facebookresearch/moco/blob/colab-notebook/colab/moco_cifar10_demo.ipynb`. Accessed: 2022-05-26.

[15] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.

[16] Reza Zanjirani Farahani and Masoud Hekmatfar. *Facility location: concepts, models, algorithms and case studies*. Springer Science & Business Media, 2009.

[17] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In *International Conference on Machine Learning*, pages 1183–1192. PMLR, 2017.

[18] Mingfei Gao, Zizhao Zhang, Guo Yu, Sercan Ö Arık, Larry S Davis, and Tomas Pfister. Consistency-based semi-supervised active learning: Towards minimizing labeling cost. In *European Conference on Computer Vision*, pages 510–526. Springer, 2020.

[19] Priya Goyal, Dhruv Mahajan, Abhinav Gupta, and Ishan Misra. Scaling and benchmarking self-supervised visual representation learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6391–6400, 2019.

[20] Guy Hacohen, Avihu Dekel, and Daphna Weinshall. Active learning on a budget: Opposite strategies suit high and low budgets. *ArXiv*, abs/2202.02794, 2022.

[21] Elad Hoffer, Tal Ben-Nun, Itay Hubara, Niv Giladi, Torsten Hoefler, and Daniel Soudry. Augment your batch: Improving generalization through instance repetition. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8126–8135, 2020.

[22] Alex Holub, Pietro Perona, and Michael C Burl. Entropy-based active learning for object recognition. In *2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, pages 1–8. IEEE, 2008.

[23] Neil Houlsby, José Miguel Hernández-Lobato, and Zoubin Ghahramani. Cold-start active learning with robust ordinal matrix factorization. In *International conference on machine learning*, pages 766–774. PMLR, 2014.

[24] Siyu Huang, Tianyang Wang, Haoyi Xiong, Jun Huan, and Dejing Dou. Semi-supervised active learning with temporal output discrepancy. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3447–3456, 2021.

[25] Shruti Jadon. Covid-19 detection from scarce chest x-ray image data using few-shot deep learning approach. In *Medical Imaging 2021: Imaging Informatics for Healthcare, Research, and Applications*, volume 11601, page 116010X. International Society for Optics and Photonics, 2021.

[26] Siddharth Karamcheti, Ranjay Krishna, Li Fei-Fei, and Christopher D Manning. Mind your outliers! investigating the negative impact of outliers on active learning for visual question answering. *arXiv preprint arXiv:2107.02331*, 2021.

[27] Jakob Nikolas Kather, Johannes Krisam, Pornpimol Charoentong, Tom Luedde, Esther Herpel, Cleo-Aron Weis, Timo Gaiser, Alexander Marx, Nektarios A. Valous, Dyke Ferber, Lina Jansen, Constantino Carlos Reyes-Aldasoro, Inka Zörnig, Dirk Jäger, Hermann Brenner, Jenny Chang-Claude, Michael Hoffmeister, and Niels Halama. Predicting survival from colorectal cancer histology slides using deep learning: A retrospective multicenter study. *PLoS Medicine*, 16, 2019.

[28] 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.

[29] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.

[30] Adrian Lang, Christoph Mayer, and Radu Timofte. Best practices in pool-based active learning for image classification. 2021.

[31] Yunfan Li, Peng Hu, Zitao Liu, Dezhong Peng, Joey Tianyi Zhou, and Xi Peng. Contrastive clustering. In *2021 AAAI Conference on Artificial Intelligence (AAAI)*, 2021.

[32] Geert Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco Ciompi, Mohsen Ghafoorian, Jeroen Awm Van Der Laak, Bram Van Ginneken, and Clara I Sánchez. A survey on deep learning in medical image analysis. *Medical image analysis*, 42:60–88, 2017.

[33] Katerina Margatina, Loic Barrault, and Nikolaos Aletras. Bayesian active learning with pretrained language models. *arXiv preprint arXiv:2104.08320*, 2021.

[34] Christoph Mayer and Radu Timofte. Adversarial sampling for active learning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3071–3079, 2020.

[35] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.

[36] Sudhanshu Mittal, Maxim Tatarchenko, Özgün Çiçek, and Thomas Brox. Parting with illusions about deep active learning. *arXiv preprint arXiv:1912.05361*, 2019.

[37] Vishwesh Nath, Dong Yang, Holger R Roth, and Daguang Xu. Warm start active learning with proxy labels and selection via semi-supervised fine-tuning. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 297–308. Springer, 2022.

[38] Kossar Pourahmadi, Parsa Nooralinejad, and Hamed Pirsiavash. A simple baseline for low-budget active learning. *arXiv preprint arXiv:2110.12033*, 2021.

[39] Tian Qiu, Guang Chen, Zi-Ke Zhang, and Tao Zhou. An item-oriented recommendation algorithm on cold-start problem. *EPL (Europhysics Letters)*, 95(5):58003, 2011.

[40] Joshua Robinson, Ching-Yao Chuang, Suvrit Sra, and Stefanie Jegelka. Contrastive learning with hard negative samples. *arXiv preprint arXiv:2010.04592*, 2020.

[41] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. *arXiv preprint arXiv:1708.00489*, 2017.

[42] Burr Settles. Active learning literature survey. 2009.

[43] 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*, pages 1308–1318. PMLR, 2020.

[44] Oriane Siméoni, Mateusz Budnik, Yannis Avrithis, and Guillaume Gravier. Rethinking deep active learning: Using unlabeled data at model training. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pages 1220–1227. IEEE, 2021.

[45] Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell. Variational adversarial active learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5972–5981, 2019.

[46] Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari S Morcos. Beyond neural scaling laws: beating power law scaling via data pruning. *arXiv preprint arXiv:2206.14486*, 2022.

[47] Jamshid Sourati, Ali Gholipour, Jennifer G Dy, Xavier Tomas-Fernandez, Sila Kurugol, and Simon K Warfield. Intelligent labeling based on fisher information for medical image segmentation using deep learning. *IEEE transactions on medical imaging*, 38(11):2642–2653, 2019.

[48] Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie, Yizhong Wang, Hannaneh Hajishirzi, Noah A Smith, and Yejin Choi. Dataset cartography: Mapping and diagnosing datasets with training dynamics. *arXiv preprint arXiv:2009.10795*, 2020.

[49] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. *arXiv preprint arXiv:2203.12602*, 2022.

[50] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herv'e J'egou. Training data-efficient image transformers & distillation through attention. In *ICML*, 2021.

[51] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, Marc Proesmans, and Luc Van Gool. Scan: Learning to classify images without labels. In *European Conference on Computer Vision*, pages 268–285. Springer, 2020.

[52] Dan Wang and Yi Shang. A new active labeling method for deep learning. In *2014 International joint conference on neural networks (IJCNN)*, pages 112–119. IEEE, 2014.

[53] Jiancheng Yang, Rui Shi, Donglai Wei, Zequan Liu, Lin Zhao, Bilian Ke, Hanspeter Pfister, and Bingbing Ni. Medmnist v2: A large-scale lightweight benchmark for 2d and 3d biomedical image classification. *arXiv preprint arXiv:2110.14795*, 2021.

[54] Ofer Yehuda, Avihu Dekel, Guy Hacohen, and Daphna Weinshall. Active learning through a covering lens. *arXiv preprint arXiv:2205.11320*, 2022.

[55] Michelle Yuan, Hsuan-Tien Lin, and Jordan Boyd-Graber. Cold-start active learning through self-supervised language modeling. *arXiv preprint arXiv:2010.09535*, 2020.

[56] Zi-Ke Zhang, Chuang Liu, Yi-Cheng Zhang, and Tao Zhou. Solving the cold-start problem in recommender systems with social tags. *EPL (Europhysics Letters)*, 92(2):28002, 2010.

[57] Evgenii Zheltonozhskii, Chaim Baskin, Alex M Bronstein, and Avi Mendelson. Self-supervised learning for large-scale unsupervised image clustering. *arXiv preprint arXiv:2008.10312*, 2020.

[58] S Kevin Zhou, Hayit Greenspan, Christos Davatzikos, James S Duncan, Bram van Ginneken, Anant Madabhushi, Jerry L Prince, Daniel Rueckert, and Ronald M Summers. A review of deep learning in medical imaging: Imaging traits, technology trends, case studies with progress highlights, and future promises. *Proceedings of the IEEE*, 2021.

[59] Zongwei Zhou. *Towards Annotation-Efficient Deep Learning for Computer-Aided Diagnosis*. PhD thesis, Arizona State University, 2021.

[60] Zongwei Zhou, Jae Shin, Ruibin Feng, R Todd Hurst, Christopher B Kendall, and Jianming Liang. Integrating active learning and transfer learning for carotid intima-media thickness video interpretation. *Journal of digital imaging*, 32(2):290–299, 2019.

[61] Zongwei Zhou, Jae Shin, Lei Zhang, Suryakanth Gurudu, Michael Gotway, and Jianming Liang. Fine-tuning convolutional neural networks for biomedical image analysis: actively and incrementally. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7340–7349, 2017.

[62] Zongwei Zhou, Jae Y Shin, Suryakanth R Gurudu, Michael B Gotway, and Jianming Liang. Active, continual fine tuning of convolutional neural networks for reducing annotation efforts. *Medical Image Analysis*, page 101997, 2021.

[63] 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. *IEEE Transactions on Knowledge and Data Engineering*, 32(4):631–644, 2019.

# A Implementation Configurations

## A.1 Data Split

PathMNIST with nine categories has 107,180 colorectal cancer tissue histopathological images extracted from Kather *et al.* [27], with 89,996/10,004/7,180 images for training/validation/testing. BloodMNIST contains 17,092 microscopic peripheral blood cell images extracted from Acevedo *et al.* [1] with eight categories, where 11,959/1,712/3,421 images for training/validation/testing. OrganAMNIST consists of the axial view abdominal CT images based on Bilic *et al.* [8], with 34,581/6,491/17,778 images of 11 categories for training/validation/testing. CIFAR-10-LT ($\rho$=100) consists of a subset of CIFAR-10 [29], with 12,406/10,000 images for training/testing.

## A.2 Training Recipe for Contrastive Learning

**Pseudocode for Our Proposed Strategy.** The algorithm 1 provides the pseudocode for our proposed hard-to-contrast initial query strategy, as elaborated in §2.

---

**Algorithm 1:** Active querying hard-to-contrast data

---

**input:**
$\mathcal{D} = \{\boldsymbol{x}_m\}_{m=1}^{M}$ {unlabeled dataset $\mathcal{D}$ contains $M$ images}
annotation budget $B$; the number of clusters $K$; batch size $N$; the number of epochs $E$
constant $\tau$; structure of encoder $f$, projection head $g$; augmentation $\mathcal{T}$
$\theta^{(e)}, e \in [1, E]$ {model parameters at epoch $e$ during contrastive learning}
**output:**
selected query $\mathcal{Q}$
$\mathcal{Q} = \varnothing$
**for** epoch $e \in \{1, \dots, E\}$ **do**
  **for** sampled minibatch $\{\boldsymbol{x}_n\}_{n=1}^{N}$ **do**
    **for all** $n \in \{1, \dots, N\}$ **do**
      draw two augmentation functions $t \sim \mathcal{T}$, $t' \sim \mathcal{T}$
      # the first augmentation
      $\tilde{\boldsymbol{x}}_{2n-1} = t(\boldsymbol{x}_n)$
      $\boldsymbol{h}_{2n-1} = f(\tilde{\boldsymbol{x}}_{2n-1})$                # representation
      $\boldsymbol{z}_{2n-1} = g(\boldsymbol{h}_{2n-1})$                # projection
      # the second augmentation
      $\tilde{\boldsymbol{x}}_{2n} = t'(\boldsymbol{x}_n)$
      $\boldsymbol{h}_{2n} = f(\tilde{\boldsymbol{x}}_{2n})$                # representation
      $\boldsymbol{z}_{2n} = g(\boldsymbol{h}_{2n})$                # projection
    **end for**
    **for all** $i \in \{1, \dots, 2N\}$ and $j \in \{1, \dots, 2N\}$ **do**
      $s_{i,j} = \boldsymbol{z}_i^\top \boldsymbol{z}_j / (\|\boldsymbol{z}_i\| \|\boldsymbol{z}_j\|)$        # pairwise similarity
      $p_{i,j} = \frac{\exp(s_{i,j})/\tau}{\sum_{n=1}^{2N} \mathbb{1}_{[n \neq i]} \exp(s_{i,n})/\tau}$        # predicted probability of contrastive pre-text task
    **end for**
    $p_{\theta^{(e)}}(y_n^*|x_n) = \frac{1}{2}[p_{2n-1,2n} + p_{2n,2n-1}]$
  **end for**
**end for**
**for** unlabeled images $\{\boldsymbol{x}_m\}_{m=1}^{M}$ **do**
  $\hat{\mu}_m = \frac{1}{E} \sum_{e=1}^{E} p_{\theta^{(e)}}(y_m^*|x_m)$
  Assign $\boldsymbol{x}_m$ to one of the clusters computed by $K$-mean($\boldsymbol{h}, K$)
**end for**
**for all** $k \in \{1, \dots, K\}$ **do**
  sort images in the cluster $K$ based on $\hat{\mu}$ in an ascending order
  query labels for top $B/K$ samples, yielding $Q_k$
  $\mathcal{Q} = \mathcal{Q} \cup \mathcal{Q}_k$
**end for**
**return** $\mathcal{Q}$

---

Table 3: Contrastive learning settings on MedMNIST and CIFAR-10-LT.

(a) MedMNIST pre-training

| config | value |
|---|---|
| backbone | ResNet-50 |
| optimizer | SGD |
| optimizer momentum | 0.9 |
| weight decay | 1e-4 |
| base learning rate[†] | 0.03 |
| learning rate schedule | cosine decay |
| warmup epochs | 5 |
| epochs | 200 |
| repeated sampling [21] | see Table 5 |
| augmentation | see Table 4 |
| batch size | 4096 |
| queue length [15] | 65536 |
| $\tau$ (equation 1) | 0.05 |

(b) CIFAR-10-LT pre-training

| config | value |
|---|---|
| backbone | ResNet-50 |
| optimizer | SGD |
| optimizer momentum | 0.9 |
| weight decay | 1e-4 |
| base learning rate[†] | 0.03 |
| learning rate schedule | cosine decay |
| warmup epochs | 5 |
| epochs | 800 |
| repeated sampling [21] | none |
| augmentation | see Table 4 |
| batch size | 512 |
| queue length [15] | 4096 |
| $\tau$ (equation 1) | 0.05 |

[†] $lr = base\_lr \times$ batchsize / 256 per the linear $lr$ scaling rule [19].

Table 4: **Data augmentations.**

(a) Augmentations for RGB images

| augmentation | value |
|---|---|
| hflip | |
| crop | [0.08, 1] |
| color jitter | [0.4, 0.4, 0.4, 0.1], p=0.8 |
| gray scale | |
| Gaussian blur | $\sigma_{min}$=0.1, $\sigma_{max}$=2.0, p=0.5 |

(b) Augmentations for OrganAMNIST

| augmentation | value |
|---|---|
| hflip | |
| crop | [0.08, 1] |
| color jitter | [0.4, 0.4, 0.4, 0.1], p=0.8 |
| rotation | degrees=45 |

**Pre-training Settings.** Our settings mostly follow [15, 14]. Table 3a summarizes our contrastive pre-training settings on MedMNIST, following [15]. Table 3a shows the corresponding pre-training settings on CIFAR-10-LT, following the official MoCo demo on CIFAR-10 [14]. The contrastive learning model is pre-trained on 2 NVIDIA RTX3090 GPUs with 24GB memory each. The total number of model parameters is 55.93 million, among which 27.97 million requires gradient backpropagation.

**Dataset Augmentation.** We apply the same augmentation as in MoCo v2 [15] on all the images of RGB modalities to reproduce the optimal augmentation pipeline proposed by the authors, including PathMNIST, BloodMNIST, CIFAR-10-LT. Because OrganAMNIST is a grey scale CT image dataset, we apply the augmentation in [3] designed for radiological images, replacing random gray scale and Gaussian blur with random rotation. Table 4 shows the details of data augmentation.

**Repeated Augmentation.** Our MoCo v2 pre-training is so fast in computation that data loading becomes a new bottleneck that dominates running time in our setup. We perform repeated augmentation on MedMNIST datasets at the level of dataset, also to enlarge augmentation space and improve generalization. [21] proposed repeated augmentation in a growing batch mode to improve generalization and convergence speed by reducing variances. This approach provokes a challenge in computing resources. Recent works [21, 50, 7] proved that fixed batch mode also boosts generalization and optimization by increasing mutiplicity of augmentations as well as parameter updates and decreasing the number of unique samples per batch, which holds the batch size fixed. Because the original contrastive learning works [13, 15] were implemented on ImageNet dataset, we attempt to simulate the quantity of ImageNet per epoch to achieve optimal performances. The details are shown in Table 5.

We only applied repeated augmentation on MedMNIST, but not CIFAR-10-LT. This is because we follow all the settings of the official CIFAR-10 demo [14] in which repeated augmentation is not employed.

Table 5: **Repeated augmentation.** For a faster model convergence, we apply repeated augmentation [21, 49, 50] on MedMNIST by reproducing the large batch size and iteration numbers.

| | # training | repeated times | # samples per epoch |
|---|---|---|---|
| ImageNet | 1,281,167 | 1 | 1,281,167 |
| PathMNIST | 89,996 | 14 | 1,259,944 |
| OrganAMNIST | 34,581 | 37 | 1,279,497 |
| BloodMNIST | 11,959 | 105 | 1,255,695 |
| CIFAR-10-LT($\rho$=100) | 12,406 | 1 | 12,406 |

## A.3 Training Recipe for MedMNIST and CIFAR-10

**Benchmark Settings.** We evaluate the initial queries by the performance of model trained on the selected initial query, and present the results in Table 1, 7 and Figure 4. The benchmark experiments are performed on NVIDIA RTX 1080 GPUs, with the following settings in Table 6.

**Cold Start Settings for Existing Active Querying Criteria.** To compare the cold start performance of active querying criteria with random selection ( Figure 1), we trained a model with the test set and applied existing active querying criteria.

Table 6: **Benchmark settings.** We apply the same settings for training MedMNIST, CIFAR-10, and CIFAR-10-LT.

| config | value |
|---|---|
| backbone | Inception-ResNet-v2 |
| optimizer | SGD |
| learning rate | 0.1 |
| learning rate schedule | reduce learning rate on plateau, factor=0.5, patience=8 |
| early stopping patience | 50 |
| max epochs | 10000 |
| augmentation | flip, p=0.5<br>rotation, p=0.5, in 90, 180, or 270 degrees<br>reverse color, p=0.1<br>fade color, p=0.1, 80% random noises + 20% original image |
| batch size | 128 |

# B   Additional Results on MedMNIST

## B.1   Label Diversity is a Significant Add-on to Most Querying Strategies

As we present in Table 1, label diversity is an important underlying criterion in designing active querying criteria. We plot the full results on all three MedMNIST datasets in Figure 6. Most existing active querying strategies became more performant and robust in the presence of label diversity.
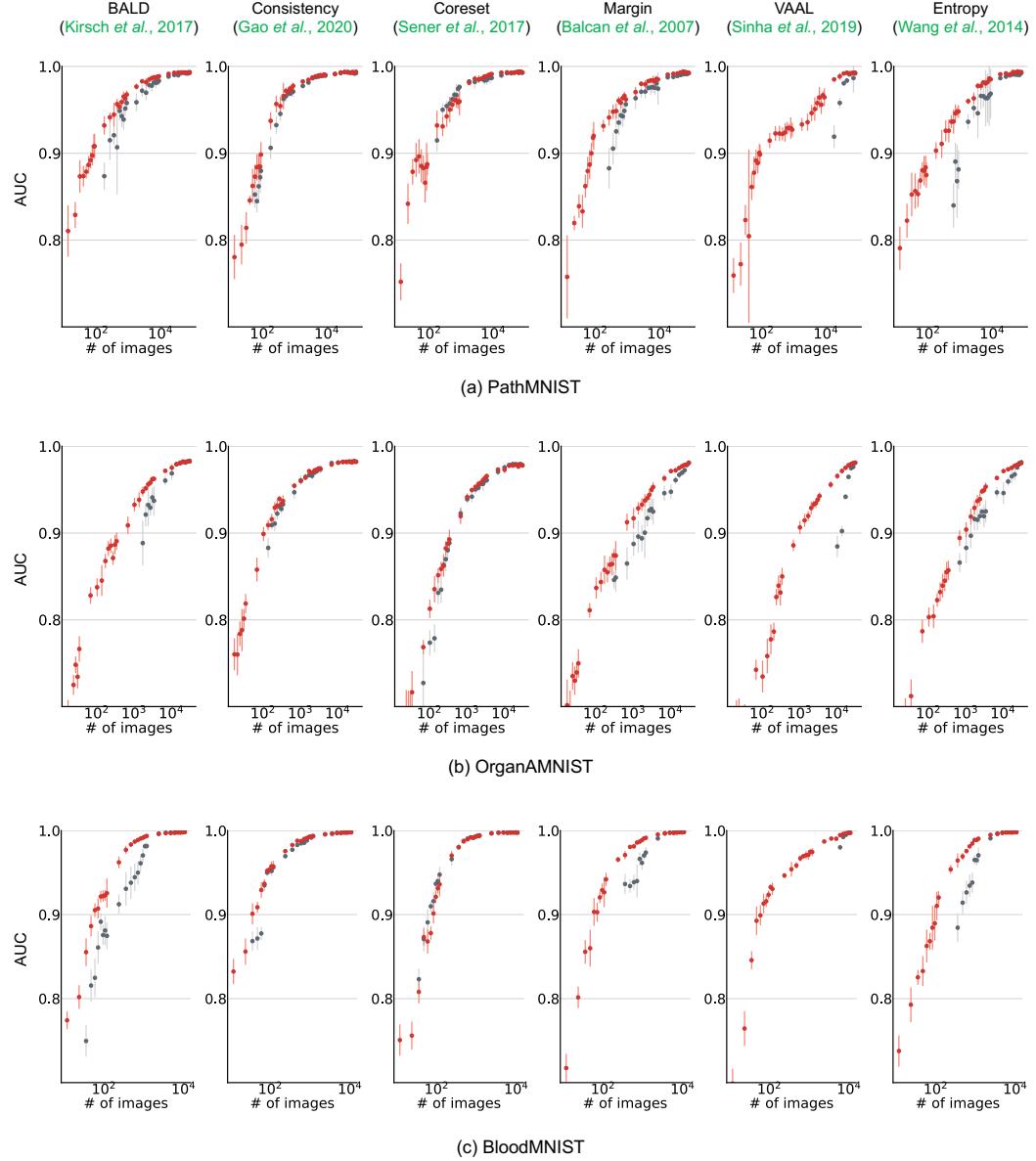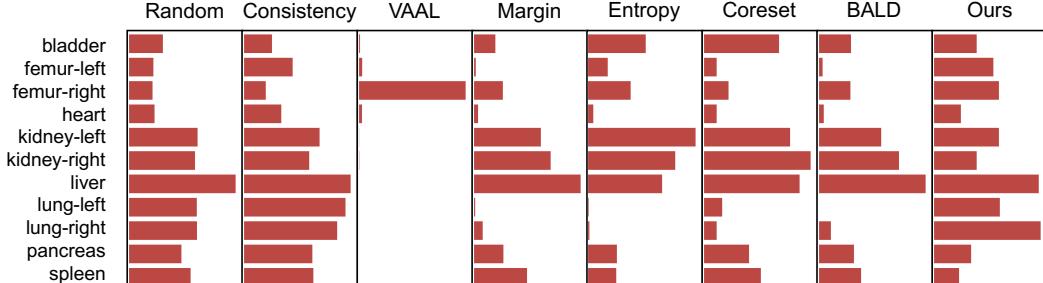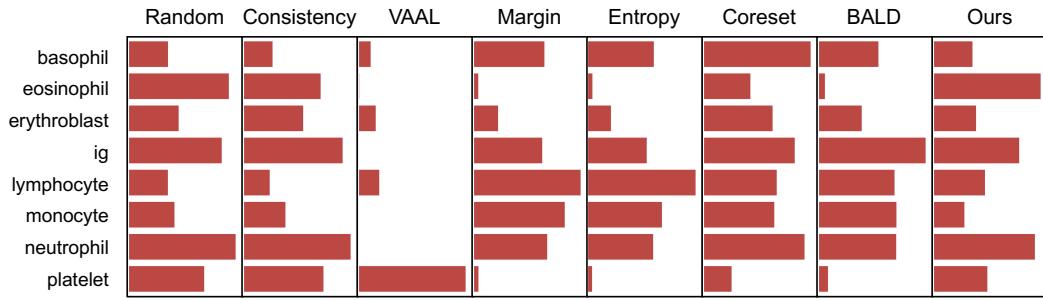


Figure 6: [Extended from Table 1] **Label diversity yields more performant and robust active querying strategies.** The experiments are conducted on three datasets in MedMNIST. The red and gray dots denote AUC scores of different active querying strategies with and without label diversity, respectively. Most existing active querying strategies became more performant and robust in the presence of label diversity, *e.g.* BALD, Margin, VAAL, and Uncertainty in particular. Some gray dots are not plotted in the low budget regime because there are classes absent in the queries due to the selection bias.

## B.2 Contrastive Features Enable Label Diversity to Mitigate Bias

Our proposed active querying strategy is capable of covering the majority of classes in most low budget scenarios by integrating K-means clustering and contrastive features, including the tail classes (*e.g.* femur-left, basophil). Compared to the existing active querying criteria, we achieve the best class coverage of selected query among at all budgets presented in Table 2.



(a) OrganAMNIST



(b) BloodMNIST

Figure 7: [Continued from Figure 2] **Our querying strategy yields better label diversity.** Random on the leftmost denotes the class distribution of randomly queried samples, which can also reflect the approximate class distribution of the entire dataset. As seen, even with a relatively larger initial query budget (691 images, 2% of OrganAMNIST, and 2,391 images, 20% of BloodMNIST), most active querying strategies are biased towards certain classes. For example in OrganAMNIST, VAAL prefers selecting data in the femur-right and platelet class, but largely ignores data in the lung, liver and monocyte classes. On the contrary, our querying strategy not only selects more data from minority classes (e.g., femur-left and basophil) while retaining the class distribution of major classes.

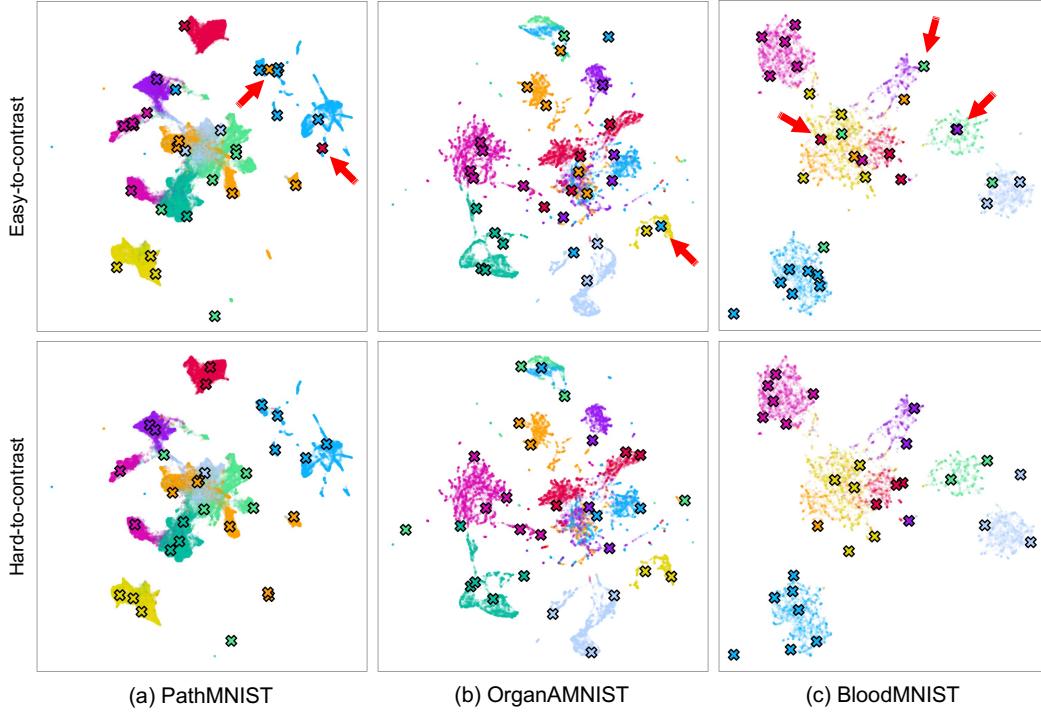Figure 8: **Visualization of $K$-means clustering and our active selection.** UMAP [35] is used to visualize the feature clustering. Colors indicate the ground truth. Contrastive features clustered by the $K$-means algorithm present a fairly clear separation in the 2D space, which helps enforce the label diversity without the need of ground truth. The crosses denote the selected easy- (top) and hard-to-contrast (bottom) data. Overall, hard-to-contrast data have a greater spread within each cluster than easy-to-contrast ones. In addition, we find that easy-to-contrast tends to select outlier classes that do not belong to the majority class in a cluster (see red arrows). This behavior will invalidate the purpose of clustering and inevitably jeopardize the label diversity.

**Selected Query Visualization.** To ease the analysis, we project the image features (extracted by a trained MoCo v2 encoder) onto a 2D space by UMAP [35]. The assigned pseudo labels have large overlap with ground truths, suggesting that the features from MoCo v2 are quite discriminative for each class. Overall, Figure 8 shows that hard-to-contrast queries have a greater spread within each cluster than easy-to-contrast ones. Both strategies can cover 100% classes. Nevertheless, we notice that easy-to-contrast selects *local outliers* in clusters: samples that do not belong to the majority class in a cluster. Such behavior will invalidate the purpose of clustering, which is to query uniformly by separating classes. Additionally, it possibly exposes the risk of introducing out-of-distribution data to the query, which undermines active learning [26].

# C  Experiments on CIFAR-10 and CIFAR-10-LT

## C.1  Label Diversity is a Significant Add-on to Most Querying Strategies

As illustrated in Table 7 and Figure 9, label diversity is an important underlying criterion in designing active querying criteria on CIFAR-10-LT, an extremely imbalanced dataset. We compare the results of CIFAR-10-LT with MedMNIST datasets Figure 6. CIFAR-10-LT is more imbalanced than MedMNIST, and the performance gain and robustness improvement of label diversity CIFAR-10-LT is significantly larger than MedMNIST. Most of the active querying strategies fail to query all the classes even at relatively larger initial query budgets.

Table 7: **Diversity is a significant add-on to most querying strategies.** AUC scores of different querying strategies are compared on CIFAR-10 and CIFAR-10-LT. In the low budget regime (*e.g.* 10% and 20% of the entire dataset), active querying strategies benefit from enforcing the label diversity of the selected data. The cells are highlighted in blue when adding diversity performs no worse than the original querying strategies. Some results are missing (marked as "-") because the querying strategy fails to sample at least one data point for each class. Results of more sampling ratios are presented in Appendix Figure 9.

| | Unif. | CIFAR-10-LT | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | | 1% (142) | 5% (710) | 10% (1420) | 20% (2841) | 30% (4261) | 40% (5682) |
| Consistency | ✓ | 78.0±1.2 | 90.0±0.1 | 91.4±1.1 | 93.4±0.2 | 93.2±0.2 | 94.6±0.2 |
| | ✗ | - | - | 67.1±17.1 | 88.6±0.3 | 90.4±0.6 | 90.7±0.2 |
| VAAL | ✓ | 80.9±1.0 | 90.3±0.5 | 92.6±0.2 | 93.7±0.4 | 93.9±0.8 | 94.5±0.2 |
| | ✗ | - | - | - | - | - | 77.3±1.6 |
| Margin | ✓ | 81.2±1.8 | 88.7±0.7 | 91.7±0.9 | 93.2±0.2 | 94.5±0.1 | 94.7±0.4 |
| | ✗ | - | - | 81.9±0.8 | 86.3±0.3 | 87.4±0.2 | 88.1±0.1 |
| Entropy | ✓ | 78.1±1.4 | 89.6±0.5 | 92.0±1.2 | 91.9±1.3 | 94.0±0.6 | 94.0±0.7 |
| | ✗ | - | 79.0±1.2 | 65.6±15.6 | 86.4±0.2 | 88.5±0.2 | 89.5±0.7 |
| Coreset | ✓ | 80.8±1.0 | 89.7±1.3 | 91.5±0.4 | 93.6±0.2 | 93.4±0.7 | 94.8±0.1 |
| | ✗ | - | - | 65.9±15.9 | 86.9±0.1 | 88.2±0.1 | 90.3±0.2 |
| BALD | ✓ | 83.3±0.6 | 90.8±0.3 | 92.8±0.1 | 90.8±2.4 | 94.0±0.8 | 94.7±0.4 |
| | ✗ | - | 76.8±2.3 | 64.9±14.9 | 84.7±0.6 | 88.0±0.5 | 88.9±0.1 |

## C.2  Contrastive Features Enable Label Diversity to Mitigate Bias

Our proposed active querying strategy is capable of covering the majority of classes in most low budget scenarios by integrating K-means clustering and contrastive features, including the tail classes (horse, ship, and truck). Compared to the existing active querying criteria, we achieve the best class coverage of selected query among at all budgets presented in Table 2. As depicted in Figure 9, our querying strategy has a more similar distribution to the overall distribution of dataset and successfully covers all the classes, with the highest proportion of minor classes (ship and truch) among random selection and all active querying methods.
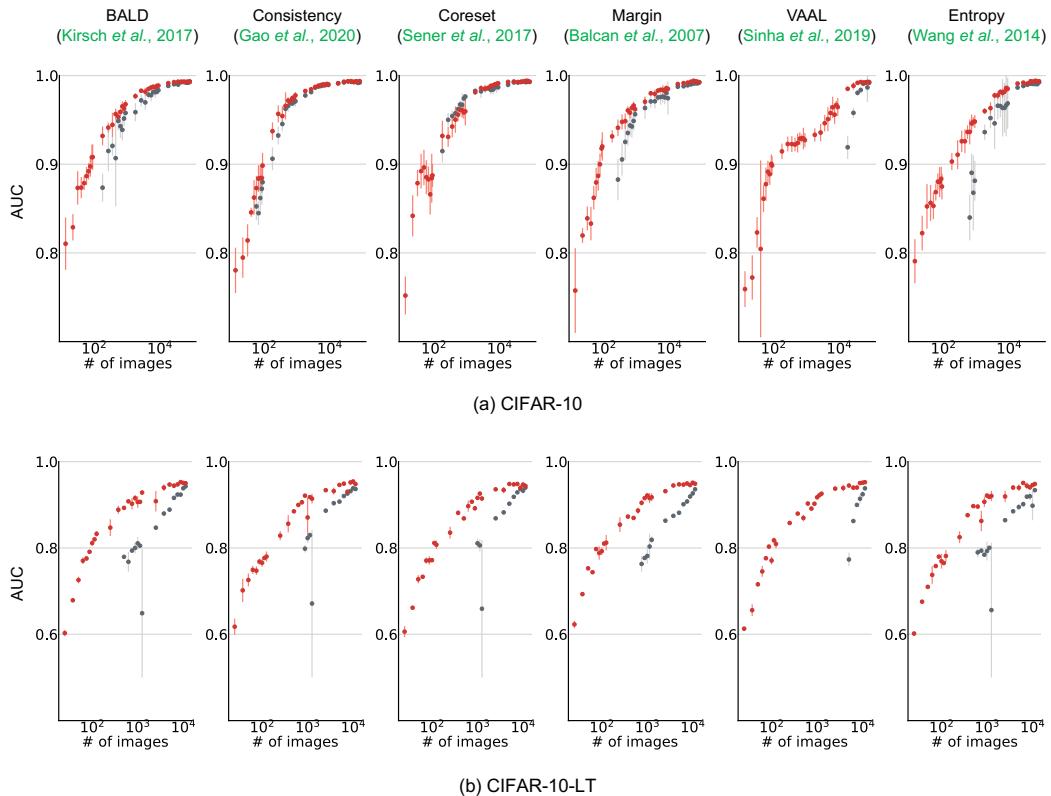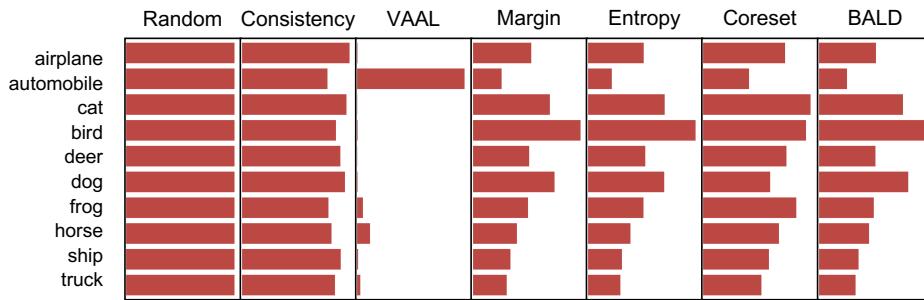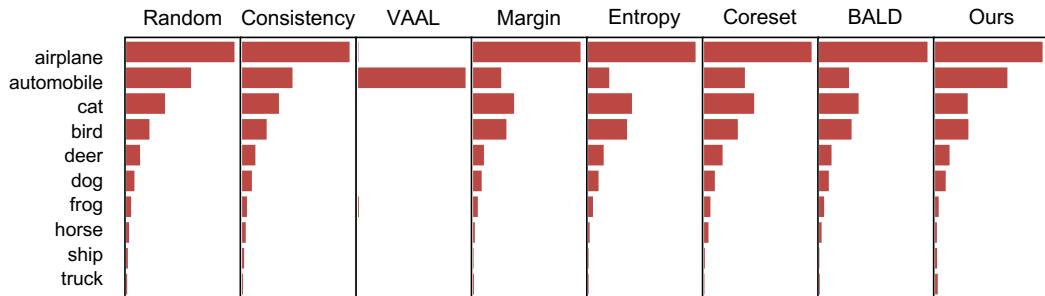
Figure 9: **Diversity yields more performant and robust active querying strategies.** The experiments are conducted on CIFAR-10-LT. The red and gray dots denote AUC scores of different active querying strategies with and without label diversity, respectively. Observations are consistent with those in medical applications (see Figure 6): Most existing active querying strategies became more performant and robust in the presence of label diversity.
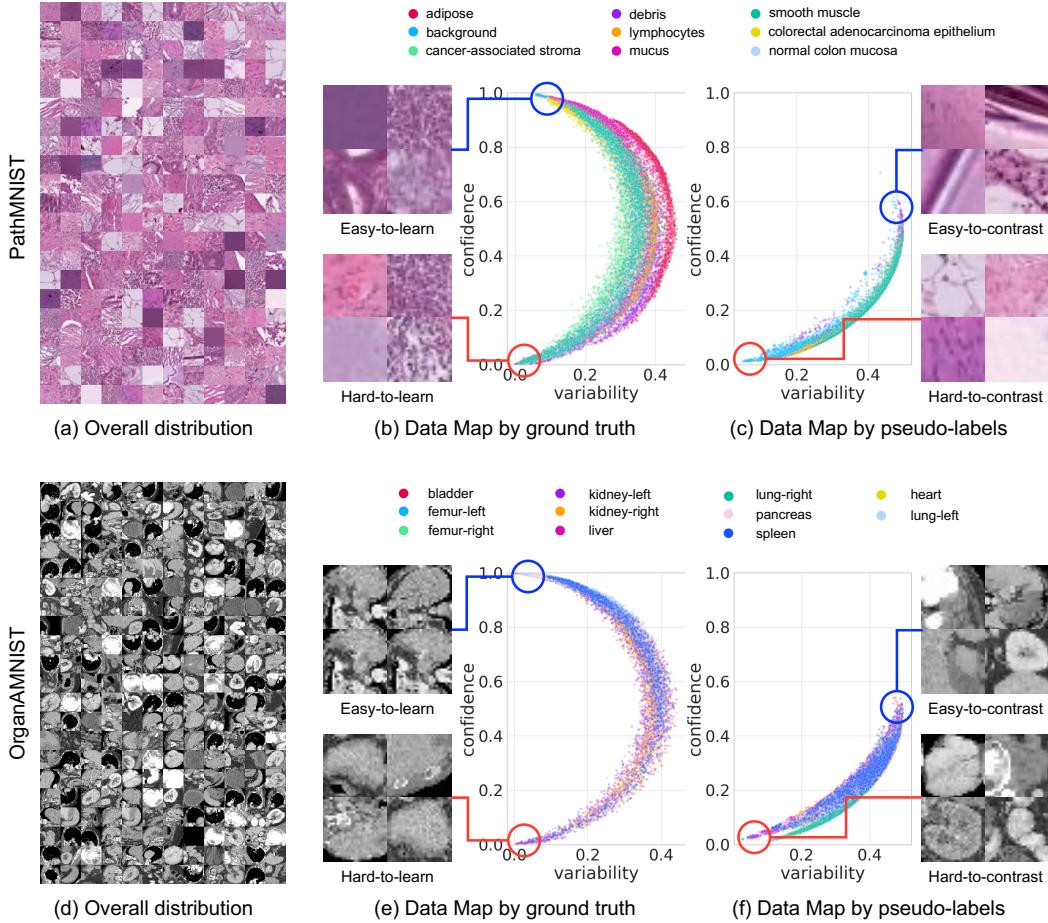
(a) CIFAR-10



(b) CIFAR-10-LT

Figure 10: **Our querying strategy yields better label diversity.** Random on the leftmost denotes the class distribution of randomly queried samples, which can also reflect the approximate class distribution of the entire dataset. As seen, even with a relatively larger initial query budget (5,000 images, 10% of CIFAR-10, and 1420 images, 10% of CIFAR-10-LT), most active querying strategies are biased towards certain classes. Our querying strategy, on the contrary, is capable of selecting more data from the minority classes such as horse, ship, and truck.

Figure 11: **Active querying based on Dataset Maps.** (a,d) PathMNIST and OrganAMNIST dataset overview. (b,e) Easy- and hard-to-learn data can be selected from the maps based on ground truths [26]. This querying strategy has two limitations: (1) requiring manual annotations and (2) data are stratified by classes in the 2D space, leading to a poor label diversity in the selected queries. (c,f) Easy- and hard-to-contrast data can be selected from the maps based on pseudo labels. This querying strategy is label-free and the selected "hard-to-contrast" data represent the most common patterns in the entire dataset. These data are more suitable for training and thus alleviate the cold start problem.
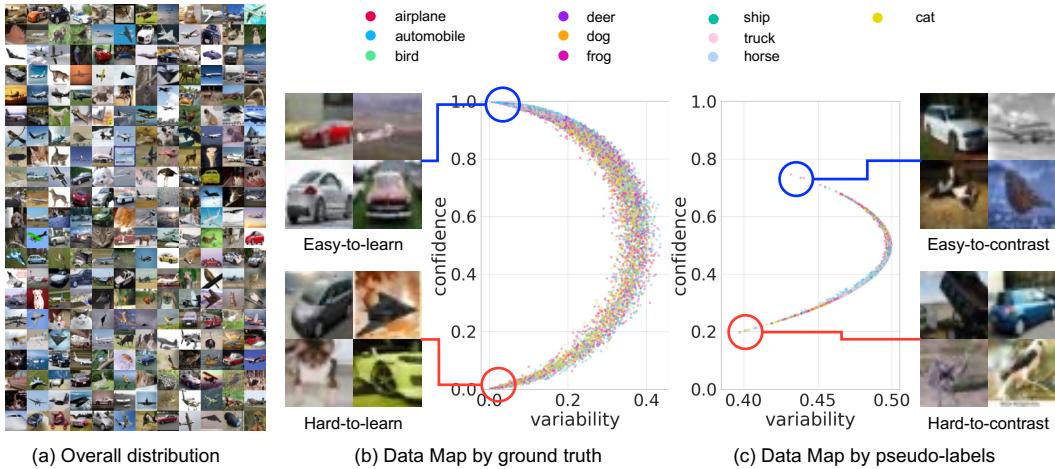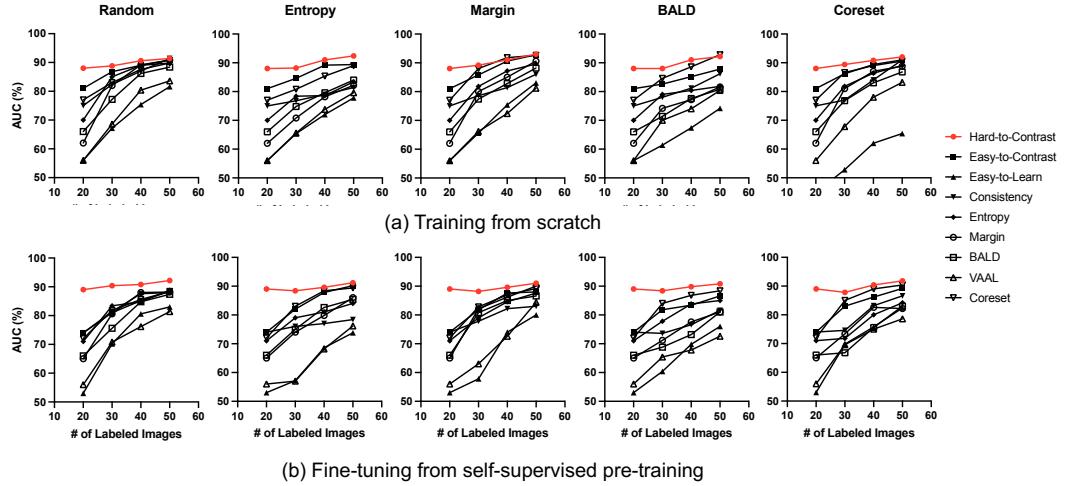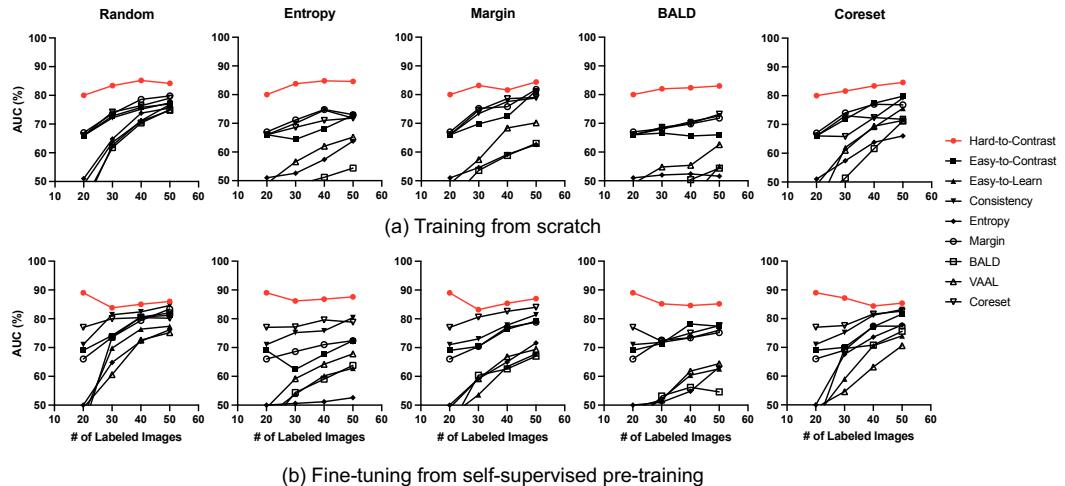
Figure 12: **Active querying based on Dataset Maps.** (a) CIFAR-10-LT dataset overview. (b) Easy- and hard-to-learn data can be selected from the maps based on ground truths [26]. This querying strategy has two limitations: (1) requiring manual annotations and (2) data are stratified by classes in the 2D space, leading to a poor label diversity in the selected queries. (c) Easy- and hard-to-contrast data can be selected from the maps based on pseudo labels. This querying strategy is label-free and the selected "hard-to-contrast" data represent the most common patterns in the entire dataset. These data are more suitable for training and thus alleviate the cold start problem.

Figure 13: **Performance of each active learning querying strategies with different initial query strategies on BloodMNIST.** Hard-to-contrast initial query strategy (red lines) outperforms other initial query strategies in every cycle of active learning. With each active learning querying strategy, the performance of the initial cycle (20 labeled images) and the last cycle (50 labeled images) are strongly correlated.



Figure 14: **Performance of each active learning querying strategies with different initial query strategies on PathMNIST.** Hard-to-contrast initial query strategy (red lines) outperforms other initial query strategies in every cycle of active learning. With each active learning querying strategy, the performance of the initial cycle (20 labeled images) and the last cycle (50 labeled images) are strongly correlated.