

---

# Towards Comparable Active Learning

---

Anonymous Authors

## Abstract

Active Learning has received significant attention in the field of machine learning. Its goal is to select the most informative samples for labeling, thereby reducing data annotation costs. However, it has been brought to attention multiple times that reported lifts from literature generalize poorly and display high variance, leading to an inconclusive landscape in Active Learning research. Based on recent insights for reliable evaluation for Active Learning, this work extends experimentation from the commonly used image domain to a wide spectrum of domains. Additionally, we provide an analysis of how many repetitions an Active Learning experiment needs in order to derive conclusive results and propose that previous benchmarks have not met the necessary number of repetitions. To the best of our knowledge, we propose the first AL benchmark that applies state-of-the-art evaluation on algorithms in 3 major domains: Tabular, Image, and Text as well as synthetic data. We report empirical results for 10 widely used algorithms on 7 real-world and 2 synthetic datasets and aggregate them into domain-specific and overall rankings of AL algorithms.

## 1 Introduction

Deep neural networks (NN) have produced state-of-the-art results on many important supervised learning tasks. Since Deep NNs usually require large amounts of labeled training data, Active Learning (AL) can be used to instead select the most informative samples out of a large pool of unlabeled data, so that only these samples need to be labeled. It has been shown that a small labeled set of this nature can be used to train well-performing models [2, 9, 15, 27].

On top of providing a principled way to label unlabeled datasets, active learning is one of the two major approaches besides semi-supervised learning to make deep learning models more data efficient by requiring only a limited set of manually labeled data. In the last decade, many different algorithms for AL have been proposed. Even though, almost every method has reported lifts over all its predecessors,<sup>1</sup> AL research faces four central difficulties: (i) The experiments are often carried out on different datasets and model architectures, hindering direct comparison, (ii) generalize poorly across different domains, (iii) the reported results can be subject to very high variance across restarts and (iv) are not always compared against important baselines like margin sampling [23]. While multiple benchmark suites have been proposed to solve (i), to the best of our knowledge, we are the first to report results on all 3 data domains of tabular, image and text. Additionally, we provide synthetic datasets to highlight principled shortcoming of existing AL algorithms. Regarding (ii) and (iii), [27] has pointed out severe inconsistencies in results of AL papers in recent years. They conducted a meta analysis of reported results of several different AL algorithms and found that all

---

<sup>1</sup>Out of all considered algorithms for this paper, only BALD [7] did not claim a new SOTA performance in their result section.

Code available at: anonymous

Table 1: Comparison of our benchmark with the existing literature. Oracle curves serve as an approximation of the best possible AL algorithm. Including the encoded versions of our datasets we reach 14 datasets. "Semi" indicates whether the paper is employing any form of self- or semi-supervised learning. A "-" for repetitions means that we could not determine how often each experiment is repeated in the respective framework.

Paper	Sampling	#Data	#Alg	Img	Txt	Tab	Synth	Semi	Oracle	Repetitions
Beck et al. [2]	batch	4	7	✓	-	-	-	-	-	-
Hu et al. [9]	batch	5	13	✓	✓	-	-	-	-	3
Zhou et al. [27]	batch	3	2	✓	✓	-	-	-	✓	5
Zhan et al. [26]	sngl+batch	35	18	-	-	✓	✓	-	✓	10-100
Munjal et al. [19]	batch	2	8	✓	-	-	-	-	-	3
Li et al. [15]	batch	5	13	✓	-	-	-	✓	-	-
Ji et al. [10]	batch	3	8	✓	-	-	-	-	-	-
Luth et al. [17]	batch	4	5	✓	-	-	-	✓	-	3
<b>Ours</b>	sngl+batch	9(14)	11	✓	✓	✓	✓	✓	✓	50

considered algorithms only provided significant lifts in their own original papers, while following literature reported performances no better than uncertainty sampling, or in some cases no better than random sampling for the same algorithm ([27] Appendix A). These outlined issues lead to an inconclusive landscape of AL algorithms, where the vast majority of reported lifts are neither statistically significant, nor prove to be generalizable. This makes it very hard to identify the best AL algorithm, or even identifying state-of-the-art algorithms. In this work we propose an evaluation protocol that was designed to handle the high variance in the performances of AL algorithms as well as being fully controllable regardless of the combination of dataset, model and AL algorithm. We base our work largely on [10], following their guidelines for a reliable evaluation of AL algorithms, while extending the number of evaluated data domains from 1 to 5.

We focus on pool-based AL where a pool of unlabeled samples is fixed at the start of each experiment and one or more samples are chosen sequentially. In addition to the default scenario of selecting a batch of samples every iteration we incorporate the single sample case into our benchmark. Batched algorithms (and benchmarks) do not have a principled advantage over single-sample AL except for speed of computation. The problem of optimizing a portfolio of unlabeled samples in each iteration is more complicated to solve and the algorithms have systematically less information per sample to work with leading to a generally worse performance, that impacts some algorithms more than others. We propose single-sample AL as an important tool to identify the best acquisition function, rather than the best combination of acquisition function and diversity regularization.

Table 1 shows a feature comparison between our proposed benchmark and several existing benchmarks in the literature. We offer our datasets in two versions - normal and encoded. An encoded dataset was pre-encoded by a self-supervised encoder model, providing a different use case for active learning. Our synthetic and text datasets (which use word embeddings in the normal setting) are exempt from this, bringing our total dataset count to 14 across 5 domains (Tabular, Image, Text, Synthetic, Encoded).

## Contributions

- C1 Study on the reproducibility of AL experiments, providing evidence that previous works might have not repeated their experiments often enough to provide reliable results
- C2 Efficient and performant algorithm for an oracle that can be constructed greedily and does not rely on search
- C3 Two novel synthetic datasets that highlight principled shortcomings of the top-performing AL algorithms
- C4 Extending the evaluation of Active Learning algorithms to 5 data domains and revealing significant differences in algorithm performance between them

## 69 2 Problem Description

70 Given two spaces  $\mathcal{X}, \mathcal{Y}$ ,  $n = l + u$  data points with  $l \in \mathbb{N}$  labeled examples  $\mathcal{L} =$   
 71  $\{(x_1, y_1), \dots, (x_l, y_l)\}$ ,  $u \in \mathbb{N}$  unlabeled examples  $\mathcal{U} = \{x_{l+1}, \dots, x_n\}$ , a model  $\hat{y} : \mathcal{X} \rightarrow \mathcal{Y}$ ,  
 72 a budget  $\mathbb{N} \ni b \leq u$  and an annotator  $A : \mathcal{X} \rightarrow \mathcal{Y}$  that can label  $x$  (we consider only hard labels  
 73 in the one-hot format). We call  $x \in \mathcal{X}$ ,  $y \in \mathcal{Y}$  predictors and labels respectively where  $(x, y)$  are  
 74 drawn from an unknown distribution  $\rho$ . Find an acquisition function  $\Omega : \mathcal{U}^{(i)}, \mathcal{L}^{(i)} \mapsto x^{(i)} \in \mathcal{U}^{(i)}$   
 75 that iteratively selects the next unlabeled point  $x^{(i)}$  for labeling

$$\begin{aligned}\mathcal{L}^{(i+1)} &\leftarrow \mathcal{L}^{(i)} \cup \{(x^{(i)}, A(x^{(i)}))\} \\ \mathcal{U}^{(i+1)} &\leftarrow \mathcal{U}^{(i)} \setminus x^{(i)}\end{aligned}$$

with  $\mathcal{U}^{(0)} = \text{seed}(\mathcal{U}, s)$  and  $\mathcal{L}^{(0)} = \{(\mathcal{U}^{(0)}, A(\mathcal{U}^{(0)}))\}$ , where  $\text{seed}(\mathcal{U}, s)$  selects  $s$  points per class for the initial labeled set.

So that the average expected loss  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  of a machine learning algorithm fitting  $\hat{y}^{(i)}$  on the respective labeled set  $\mathcal{L}^{(i)}$  is minimal:

$$\min \frac{1}{B} \sum_{i=0}^B \mathbb{E}_{(x,y) \sim \rho} \ell(y, \hat{y}^{(i)})$$

## 76 3 Related Work

77 Multiple benchmark suites have already been proposed for Active Learning: The authors of [2], [19],  
 78 [15], [10] and [17] focus exclusively on batch AL in the image domain. While [2] discuss a new  
 79 metric to measure AL performance, which they call “Label Efficiency” and provide experiments  
 80 on many common configurations of data preparation, model training, and other hyperparameters,  
 81 [15] focuses on combined approaches of AL and semi-supervised learning. The authors of [9] study  
 82 models that are trained on actively learned datasets in the image and text domain. They test for  
 83 several different properties of the models including robustness, response to compression techniques  
 84 and final performance. [27] proposed an oracle algorithm for AL that uses Simulated Annealing  
 85 search to approximate a solution for the optimal subset of labeled data. Additionally, they study  
 86 the generalization behavior of subsets of labeled data in the text an image domain. The two closest  
 87 related works to this benchmark are [10] and [17], who also discussed the problems of evaluating  
 88 AL algorithms under many forms of variance. We largely adapt the proposed guidelines of [10] and  
 89 extend their work to multiple domains, batch sizes and comparisons. This work also pays attention  
 90 to the so-called “pitfalls” of AL evaluation mentioned in [17] with the exception of P1 - Controlling  
 91 the data distribution of the tested datasets. We leave this aspect for future versions of this benchmark.  
 92 Two preceding benchmarks have also proposed an oracle algorithm to find a near-optimal labeled  
 93 set  $\mathcal{L}$ , however, both algorithms rely on search and are for that reason considerably slower than our  
 94 proposed greedy oracle.

## 95 4 Methodology

### 96 4.1 Evaluation Protocol

97 Following [27], the quality of an AL algorithm is evaluated by an “anytime protocol” that incorpo-  
 98 rates classification performance at every iteration, as opposed to evaluating final performance after  
 99 the budget is exhausted. We employ the normalized area under the accuracy curve (AUC):

$$\text{AUC}(\mathcal{D}_{\text{test}}, \hat{y}, B) := \frac{1}{B} \sum_{i=1}^B \text{Acc}(\mathcal{D}_{\text{test}}, \hat{y}^{(i)}) \quad (1)$$

100 The AUC incorporates performance in early stages (low budget) as well as capabilities to push the  
 101 classifier in later stages (high budget). AL algorithms have to perform well in both scenarios.

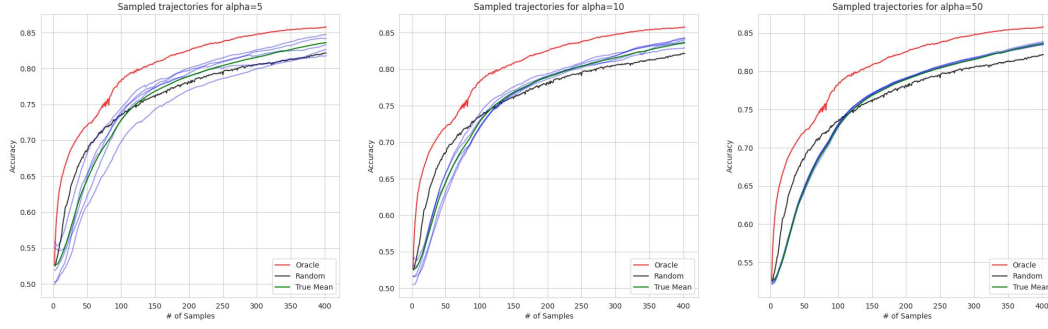


Figure 1: Random draws from a pool of 100 runs for margin sampling on the Splice dataset with different numbers of repetitions ( $\alpha = \{5, 10, 50\}$ ). Green curves are the mean performance of all 100 runs, while the samples are blue. Even with 5 or 10 repetitions, we can observe that single draws for margin sampling display below-random performance (black), while the true mean should be above random.

Since AUC is still influenced by the budget, we define a set of rules to set this hyperparameter upfront, so that we are not favoring a subset of algorithms by handcrafting a budget. In this work, we choose the budget per dataset to be the first point at which one of 2 stopping conditions apply: (i) an algorithm (except Oracle) manages to reach 99% of the full-dataset-performance (using the smallest query size) or (ii) the best algorithm (except oracle) did not improve the classifier’s accuracy by at least 2% in the last 20% of iterations.

To mimic the leave-one-out protocol for cross-validation, we will restart each experiment multiple times. Each restart retains the train/test split (often given by the dataset itself), but introduces a new validation split that is sampled from the entire dataset (not just  $\mathcal{L}$ ).

#### 4.2 Why we need 50 restarts

To evaluate how many restarts are necessary to obtain conclusive and reproducible results in an AL experiment, we set the number of restarts for margin sampling on the Splice dataset to 100. This allows us firstly, to obtain a very strong estimation of the “true” average performance of margin sampling on this dataset and secondly, to draw subsets from this pool of 100 runs. Setting the size of our draws to  $\alpha$  and sampling uniformly, we can approximate a cross-validation process with  $\alpha$  restarts. Each of these draws can be interpreted as a reported result in AL literature where the authors employed  $\alpha$  restarts. Figure 1 shows the “true” mean performance of margin sampling (green) in relation to random sampling (black) and the oracle performance (red). We display 5 random draws of size  $\alpha$  in blue. We can observe that even for a relatively high number of restarts the variance between the samples is extremely high, resulting in some performance curves being worse than random and some being significantly better. When setting  $\alpha = 50$  we observe all samples to converge close to the true mean performance. In addition to this motivating example, we carried out our main evaluation (Tab. 3) multiple times by uniformly sampling 3 random from our 50 available runs and comparing the results. We found significant differences in the performance of acquisition functions on individual datasets, as well as permutations in the final ranking. This partly explains the ongoing difficulties in reproducing results for AL experiments and benchmarks. This details can be found in App. B. For this benchmark we always employ 50 restarts of the same experiment.

#### 4.3 Seeding vs. Restarts

Considering the high computational cost of 50 repetitions, another approach to ensure reproducibility would be to reduce the amount of variance in the experiment by keeping as many subsystems (weight initialization, data splits, etc.) as possible fixed with specialized seeding.

We describe a novel seeding strategy in Appendix E that creates 3 separate Random Number Generators (RNG) based on 3 different seeds. In short, we introduce three different seeds:  $s_\Omega$  for the AL algorithm,  $s_{\mathcal{D}}$  for dataset splitting and mini-batch sampling, and  $s_\theta$  for model initialization and sampling of dropout masks. Unless stated otherwise, we will keep  $s_\Omega$  fixed, while  $s_{\mathcal{D}}$  and  $s_\theta$  are incremented by 1 between restarts to introduce stochasticity into our framework. While this seeding strategy is capable of controlling the amount variance in the experiment, previous works have noted that an actively sampled, labeled set does not generalize well between model architectures or even different initializations of the same model ([27, 16]), reducing its value in practice and providing a bad approximation of the quality of an AL algorithm. Hence, we opt for letting the subsystems vary (by increasing  $s_{\mathcal{D}}$  and  $s_\theta$ ) and combine that with a high number of restarts to obtain a good average of the generalization performance of each AL algorithm. Where a high number of restarts is computationally not feasible, we advise to additionally keep either  $s_{\mathcal{D}}$  or  $s_\theta$  (or both) fixed.

#### 4.4 Datasets

A detailed description of the preprocessing of each dataset can be found in Appendix I.

**Tabular:** AL research conducted on tabular data is sparse (only [1] from the considered baseline papers). We, therefore, introduce a set of tabular datasets that we selected according to the following criteria: (i) They should be solvable by medium-sized models in under 1000 samples, (ii) the gap between most AL algorithms and random sampling should be significant (potential for AL is present) and (iii) the gap between the AL algorithms and our oracle should also be significant (research on these datasets can produce further lifts). We use **Splice**, **DNA** and **USPS** from LibSVMTools [20].

**Image:** We use **FashionMNIST** [24] and **Cifar10** [13], since both are widely used in AL literature.

**Text:** We use **News Category** [18] and **TopV2** [6]. Text datasets have seen less attention in AL research, but most of the papers that evaluate on text ([9], [27]) use at least one of these datasets.

We would like to point out that these datasets are selected for speed of computation (both in terms of number of features and necessary budget to solve the dataset). However, similar to our argumentation for picking smaller classifiers, we are solely focused on comparing different AL algorithms in this paper and do not aim to develop novel classification models on these datasets. Our assumption is that a well-performing algorithm in our benchmark will also generalize well to larger real-world datasets, because we included multiple different data domains and classifier sizes in our experiments.

Adapting the experimental setting from [8], we offer all our datasets in the un-encoded (normal) setting as well as pre-encoded by a fixed embedding model that was trained by unsupervised contrastive learning. The text datasets are an exception to this, as they are only offered in their encoded form. The pre-encoded datasets enable us to test single-sample algorithms on more complex datasets like **Cifar10** and **FashionMnist**. They also serve the purpose of investigating the interplay between self-supervised learning techniques and AL, as well as alleviating the cold-start problem described in [17] as they require a way smaller seed set. The classification model for every encoded dataset is a single linear layer with softmax activation. The embedding model was trained with the SimCLR [5] algorithm adopting the protocol from [8]. To ensure that enough information from the data is encoded by our embedding model, the quality of embeddings during pretext training was measured after each epoch. We attached a linear classification head to the encoder, fine-tuned it to the data and evaluated this classifier for test accuracy, mirroring our AL setup for embedded datasets. The checkpoint of each encoder model will be provided together with the framework.

**Synthetic Data** Existing AL algorithms can broadly be categorized into two types, uncertainty [23, 7] - and geometric-approaches [22, 8, 1]. Both types have principled shortcomings in terms of the utilized information that makes them unsuitable for certain data distributions. To test for these shortcomings, we created two synthetic datasets, namely "Honeypot" and "Diverging Sine", that are hard to solve for methods focused on the classifier's decision boundary or data

184 clustering respectively. To avoid algorithms memorizing these datasets they are generated during  
 185 the experiment, depending on  $s_{\mathcal{D}}$ .

186

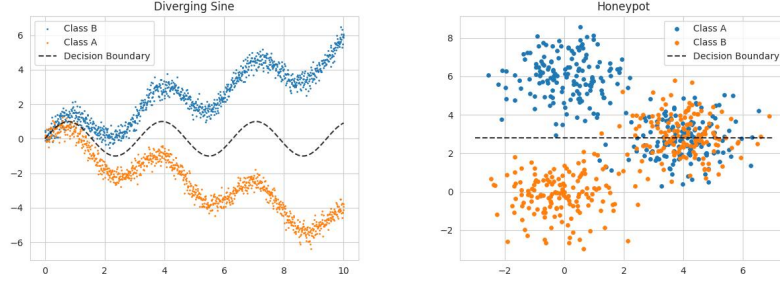


Figure 2: Synthetic "Honeypot" and "Diverging Sine" datasets. The decision boundary is not part of the dataset and serves only as a visual guide.

187 Honeypot creates to two easy to distinguish clusters with 150 samples each and one overlap-  
 188 ping "honeypot" that represents a noisy region of the dataset with potentially miss-labeled, miss-  
 189 measured or generally adverse samples. This honeypot contains 150 samples of each class, creating  
 190 a balance of 50% beneficial samples and 50% adverse samples in the dataset. The honeypot is  
 191 located on the likely decision boundary of a classifier that is trained on the beneficial samples to  
 192 maximize its negative impact on purely uncertainty based acquisition functions. Diverging Sine  
 193 samples the datapoints for each class from two diverging sinusoidal functions that are originating  
 194 from the same y-intercept. This creates a challenging region one the left hand side, where a lot of  
 195 datapoints need to be sampled and an easy region on the right hand side, where very few datapoints  
 196 are enough. The repeating nature of a sin function encourages diversity based acquisition functions  
 197 to equally sample the entire length, drastically oversampling the right hand side of the dataset. Each  
 198 class has 500 datapoints.

199 Every dataset has a fixed size for the seed set of 1 sample per class, with the only exceptions be-  
 200 ing un-encoded FashionMnist and Cifar10 with 100 examples per class to alleviate the cold-start  
 201 problem in these complex domains.

## 202 4.5 Classification Model

203 The model  $\hat{y}$  can be trained in  
 204 two ways. Either the param-  
 205 eters of the model are reset to  
 206 a fixed initial setting  $\hat{y}^{(0)}$  after  
 207 each AL iteration and the classi-  
 208 fier is trained from scratch with  
 209 the updated labeled set  $\mathcal{L}^{(i)}$ , or  
 210 the previous state  $\hat{y}^{(i-1)}$  is re-  
 211 tained and the classifier is fine-  
 212 tuned on  $\mathcal{L}^{(i)}$  for a reduced num-  
 213 ber of epochs. In this work,  
 214 we use the fine-tuning method  
 215 for un-encoded datasets to save  
 216 computational time, while we  
 217 use the from-scratch training for

Table 2: Overview of available batch sizes for each dataset

	B	1	5	20	50	100	500	1000
Enc. DNA	40	o	o					
Honeypot	60	o	o					
Diverging Sine	60	o	o					
Enc. Splice	100	o	o	o	o			
TopV2	200	o	o	o	o			
Splice	400	o	o	o	o	o		
DNA	300	o	o	o	o	o		
USPS	400	o	o	o	o	o		
Enc. Cifar10	450	o	o	o	o	o		
Enc. FMnist	500	o	o	o	o	o		
Enc. USPS	600	o	o	o	o	o		
News	3K			o	o	o	o	
FMnist	10K						o	o
Cifar10	10K						o	o

218 embedded datasets since they have very small classifiers and this approach generally produces better  
 219 results. Our fine-tuning scheme always trains for at least one epoch and employs an aggressive early  
 220 stopping with a patience of 2 afterwards.

## 4.6 Batch Sizes

We selected batch sizes for each dataset to accommodate the widest range possible that results in a reasonable runtime for low batch sizes and allows for at least 4 round of data acquisition for high batch sizes. The available batch sizes per dataset can be found in Table 2.

## 4.7 Realism vs. Variance

We would like to point out that some design choices for this framework prohibit direct transfer of our results to practical applications. This is a conscious choice, as we think that this is a necessary trade-off between realism and experiment variance. We would like to highlight the following design decisions:

- (i) Creating test and validation splits from the full dataset rather than only the labeled seed set. Fully fledged test and validation splits are unobtainable in practice, but they provide not only a better approximation of algorithm performance, but also a better foundation for hyperparameter tuning, which is bound to reduce variance in the experiment.
- (ii) Choosing smaller classifiers instead of SOTA models. Since we are not interested in archiving a new SOTA in any classification problem, we instead opt to use smaller classifiers for the following reasons: Smaller classifiers generally exhibit more stable training behavior, on average require fewer sampled datapoints to reach their full-dataset-performance and have faster training times. For every dataset, the chosen architecture’s hyperparameters are optimized to archive maximum full-dataset performance. Generally, we use MLPs for tabular, RestNet18 for image and BiLSTMs for text datasets. Every encoded dataset is classified by a single linear layer with softmax activation. For an overview of architectures and hyperparameters please refer to Appendix I.

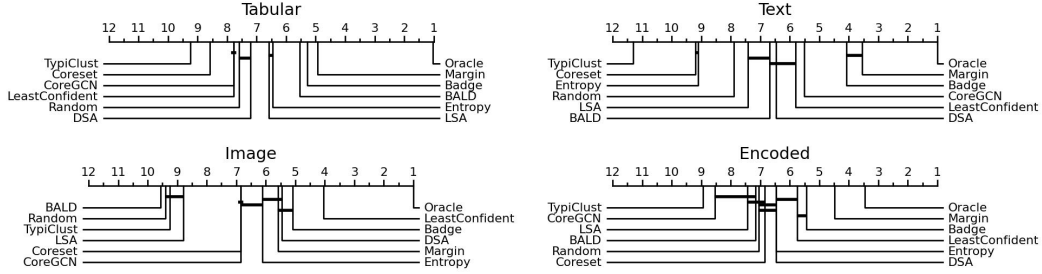
## 4.8 Greedy Oracle Algorithm

Posing Active Learning as a combinatorial problem, the oracle set  $\mathcal{O}_b$  for a given dataset, model, and training procedure is the set that induces the highest AUC score for a given budget. However, since this problem is computationally infeasible for realistic datasets, previous works have proposed approximations to this oracle sequence. [27] used simulated annealing to search for the optimal subset and used the best solution found after a fixed time budget. Even though their reported performance curves display a significant lift over all other acquisition functions, we found the computational cost of reproducing this oracle for all our datasets to be prohibitive (The authors reported the search to take several days per dataset on 8 V100 GPUs). In this paper, we propose a greedy oracle algorithm that constructs an approximation of the optimal set in an iterative fashion. Our oracle algorithm evaluates every data point  $u_k = \text{unif}(\mathcal{U})$   $k \in [1 \dots \tau]$  in a subsample of unlabeled points by fitting the classifier  $\hat{y}$  on  $\mathcal{L}^{(i)} \cup \{u_k\}$  and directly measuring the resulting test performance. The data point with the best test performance is selected and added to the labeled pool for that iteration. We noticed that this oracle is over-specializing on the test set, resulting in stagnating or even decreasing performance curves in later AL iterations. This can happen, for example, if the oracle picked a labeled set that enables the classifier to correctly classify a big portion of easy samples in the test set, but now fails to find the next **single** unlabeled point that would enable the classifier to succeed on one of the hard samples. This leads to a situation, where no point can immediately incur an increase in test performance and therefore the selected data point can be considered random. To circumvent this problem, we use margin sampling [23] as a fallback option for the oracle. Whenever the oracle does not find an unlabeled point that results in an increase in performance, it defaults to margin sampling in that iteration. The resulting greedy algorithm constructs an approximation of the optimal labeled set that consistently outperforms all other algorithms by a significant margin, while requiring relatively low computational cost ( $\mathcal{O}(B\tau)$ ). We fix  $\tau = 20$  in this work, as this gave us already a significant lift and we expect diminishing returns for larger  $\tau$ . The pseudocode for our oracle can be found in App. J. Even though our proposed algorithm is more efficient than other approaches, the computational costs for high budget datasets like Cifar10 and FashionMnist meant that we could not compute the oracle for all 10000 datapoints. To still provide an oracle for these two datasets, we select two points per iteration instead of one and stop the oracle computation at a budget of 5000.

Table 3: Performances for acquisition functions on real-world datasets, aggregated for un-encoded and encoded datasets. Performance is shown as average ranks over restarts (1.0 is the best rank). Algorithms are sorted by aggregated performance on un-encoded datasets.

	Splice	DNA	USPS	Cfr10	FMnist	TopV2	News	Un-enc.	Enc.
Oracle	1.0 $\pm$ 0.01	1.0 $\pm$ 0.01	1.0 $\pm$ 0.0	1.0 $\pm$ 0.0	1.0 $\pm$ 0.0	1.0 $\pm$ 0.01	1.0 $\pm$ 0.0	1.0	2.0
Margin	6.6 $\pm$ 0.02	4.3 $\pm$ 0.01	2.1 $\pm$ 0.01	6.3 $\pm$ 0.01	4.4 $\pm$ 0.0	2.4 $\pm$ 0.01	3.7 $\pm$ 0.0	4.3	4.2
Badge	5.2 $\pm$ 0.01	6.3 $\pm$ 0.01	2.9 $\pm$ 0.01	5.2 $\pm$ 0.01	4.7 $\pm$ 0.0	3.3 $\pm$ 0.01	3.5 $\pm$ 0.0	4.5	5.4
LeastConf	9.2 $\pm$ 0.02	10.3 $\pm$ 0.02	8.1 $\pm$ 0.02	2.1 $\pm$ 0.01	4.0 $\pm$ 0.0	7.9 $\pm$ 0.02	3.0 $\pm$ 0.01	6.4	6.5
DSA	7.4 $\pm$ 0.02	7.3 $\pm$ 0.01	7.5 $\pm$ 0.01	5.4 $\pm$ 0.01	5.1 $\pm$ 0.0	6.0 $\pm$ 0.02	7.3 $\pm$ 0.01	6.6	6.7
BALD	4.0 $\pm$ 0.01	4.7 $\pm$ 0.01	5.4 $\pm$ 0.01	12.0 $\pm$ 0.01	7.6 $\pm$ 0.0	7.6 $\pm$ 0.02	5.0 $\pm$ 0.0	6.6	7.6
CoreGCN	6.9 $\pm$ 0.01	4.9 $\pm$ 0.01	10.4 $\pm$ 0.01	7.6 $\pm$ 0.01	6.5 $\pm$ 0.01	4.0 $\pm$ 0.01	6.8 $\pm$ 0.0	6.7	8.2
Entropy	6.6 $\pm$ 0.02	3.9 $\pm$ 0.01	7.6 $\pm$ 0.01	7.6 $\pm$ 0.01	4.9 $\pm$ 0.01	9.8 $\pm$ 0.02	9.6 $\pm$ 0.0	7.1	6.5
LSA	6.1 $\pm$ 0.01	6.8 $\pm$ 0.01	5.3 $\pm$ 0.01	7.7 $\pm$ 0.01	10.6 $\pm$ 0.01	7.5 $\pm$ 0.01	7.3 $\pm$ 0.01	7.3	7.5
Random	9.0 $\pm$ 0.01	9.3 $\pm$ 0.01	5.3 $\pm$ 0.01	8.4 $\pm$ 0.01	11.1 $\pm$ 0.0	7.9 $\pm$ 0.01	8.0 $\pm$ 0.0	8.4	6.9
Coreset	7.1 $\pm$ 0.01	9.0 $\pm$ 0.01	10.5 $\pm$ 0.01	6.8 $\pm$ 0.01	7.1 $\pm$ 0.0	8.5 $\pm$ 0.02	10.8 $\pm$ 0.01	8.5	7.2
TypiClust	8.8 $\pm$ 0.01	10.2 $\pm$ 0.01	12.0 $\pm$ 0.02	7.9 $\pm$ 0.01	11.0 $\pm$ 0.01	12.0 $\pm$ 0.02	12.0 $\pm$ 0.01	10.5	9.2

Figure 3: Ranks of each acquisition function aggregated by domain. Horizontal bars indicate a **non**-significant rank difference. The significance is tested via a paired-t-test with  $\alpha = 0.05$ .



271 The rest of the curve is forecast with a simple linear regression that asymptotically approaches the  
272 upper bound performance of the dataset. A detailed description can be found in App. G.

## 273 5 Experiments

### 274 5.1 Implementation Details

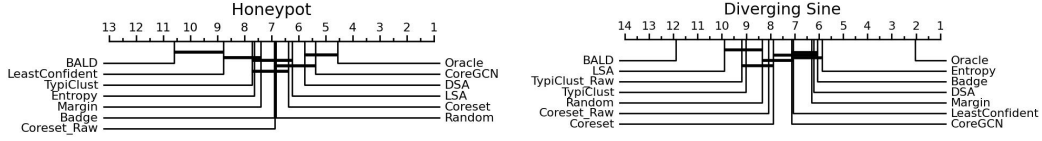
275 At each iteration  $i$  the acquisition function  $\Omega$  picks an unlabeled datapoint based on a fixed set of  
276 information  $\{\mathcal{L}^{(i)}, \mathcal{U}^{(i)}, B, |\mathcal{L}^{(i)}| - |\mathcal{L}^{(1)}|, \text{acc}^{(i)}, \text{acc}^{(1)}, \hat{y}^{(i)}, \text{opt}_{\hat{y}}\}$ , where  $\text{opt}_{\hat{y}}$  is the optimizer  
277 used to fit  $\hat{y}^{(i)}$ . This set grants full access to the labeled and unlabeled set, as well as all parameters  
278 of the classifier and the optimizer. Additionally, we provide meta-information, like the size of the  
279 seed set through  $|\mathcal{L}^{(i)}| - |\mathcal{L}^{(1)}|$ , the remaining budget though the addition of  $B$  and the classifiers  
280 potential though  $\text{acc}^{(1)}$  and  $\text{acc}^{(i)}$ . We allow acquisition functions to derive information from this  
281 set, e.g. predictions of the classifier  $\hat{y}^{(i)}(x)$ ;  $x \in \mathcal{U}^{(i)} \cup \mathcal{L}^{(i)}$ , clustering, or even training additional  
282 models. However, the algorithm may not incorporate external information e.g. other datasets,  
283 queries to recover additional labels, additional training steps for  $\hat{y}$ , or the test/validation set.  
284 For our study we selected acquisition functions with good performances reported by multiple  
285 different sources that can work with the set of information stated above. For a list of all acquisition  
286 functions, please refer to Table 3, with detailed descriptions being found in Appendix A.

### 288 5.2 Results

289 Storing the results of every run in a separate file allows us to create various aggregates. Apart from  
290 plotting standard performance curves and reporting their AUC values per dataset (Fig. 4), we pri-  
291 marily rely on ranks to aggregate the performance of an acquisition function across datasets. For  
292 each dataset, the AUC values of all acquisition functions are sorted and assigned a rank based on  
293 position, with the best rank being 1. These ranks can safely be averages across datasets as they are



Figure 4: Results for all acquisition functions on both synthetic datasets.



no longer subjected to scaling differences of each dataset. In Table 3 we provide the rank of each acquisition function per dataset and averaged for each (un-)encoded dataset. Please note, that for Tab 3 we are averaging not only over runs, but also over query sizes per dataset. For the results per query size, please refer to App. C.

Additionally, we provide aggregated ranks per data domain in Fig. 3, revealing varying performances of acquisition functions like DSA/LSA, BALD and CoreGCN. The leading algorithms Margin-Sampling and BADGE show very stable performance across all tested data domains.

Results for the HoneyPot dataset reveal expected shortcomings of uncertainty sampling algorithms like margin, entropy and least confident sampling as well as BALD. In addition, BADGE is underperforming for this dataset compared to real-world data. Results for Diverging Sine also confirm expected behavior, as clustering algorithms (Coreset, TypiClust) fall behind uncertainty algorithms (Entropy-, Margin-Sampling), with the exception of BADGE.

We provide a very small ablation study on the importance of the embeddings by testing a version of Coreset and TypiClust on this dataset that does not use the embeddings produced by the classification model, but rather clusters the data directly. "Coreset Raw" and "TypiClust Raw" both perform worse than their embedding-based counterpart.

## 6 Conclusion

As stated in contribution C4, our results on real-world data shows significant differences in the performance of the tested algorithms between data domains. Not only do some algorithms overperform on some domains (like least confidence sampling on Images), but the Top-3 of algorithms (except Oracle) does not contain the same three algorithms for any two domains. This highlights the dire need for diverse data domains in AL benchmarking. Furthermore, our results on HoneyPot reveal principled shortcomings for the two best algorithms in BADGE and margin sampling. Both are vulnerable to adverse samples or simply measurement noise, which highlights the need for further research in this area.

Finally, the fact that BADGE is able to perform well on Diverging Sine highlights the importance of embeddings for the clustering algorithms, as the so-called gradient embedding from BADGE seems to be able to encode uncertainty information, guiding the selection into the left hand regions of the dataset. We also show that embeddings are generally useful for this dataset, by providing results for "Coreset Raw" and "TypiClust Raw".

In conclusion, we strongly advocate to test newly proposed AL algorithms not only on a wide variety of real data domains, but also to pay close attention to the HoneyPot and Diverging Sine datasets to reveal principled shortcomings of the algorithm in question.

327 **Acknowledgement** anonymous

## 328 **References**

- 329 [1] Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal.  
330 Deep batch active learning by diverse, uncertain gradient lower bounds. In *International Con-*  
331 *ference on Learning Representations*, 2020.
- 332 [2] Nathan Beck, Durga Sivasubramanian, Apurva Dani, Ganesh Ramakrishnan, and Rishabh  
333 Iyer. Effective evaluation of deep active learning on image classification tasks. *arXiv preprint*  
334 *arXiv:2106.15324*, 2021.
- 335 [3] Razvan Caramalau, Binod Bhattacharai, and Tae-Kyun Kim. Sequential graph convolutional net-  
336 work for active learning. In *Proceedings of the IEEE/CVF conference on computer vision and*  
337 *pattern recognition*, pages 9583–9592, 2021.
- 338 [4] Akshay L Chandra. Deep active learning toolkit for image classification in pytorch. <https://github.com/ac121/deep-active-learning-pytorch>, 2021.
- 340 [5] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework  
341 for contrastive learning of visual representations. In *International conference on machine*  
342 *learning*, pages 1597–1607. PMLR, 2020.
- 343 [6] Xilun Chen, Asish Ghoshal, Yashar Mehdad, Luke Zettlemoyer, and Sonal Gupta. Low-  
344 resource domain adaptation for compositional task-oriented semantic parsing. In *Proceedings*  
345 *of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.  
346 Association for Computational Linguistics, 2020.
- 347 [7] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image  
348 data. In *International conference on machine learning*, pages 1183–1192. PMLR, 2017.
- 349 [8] Guy Hacohen, Avihu Dekel, and Daphna Weinshall. Active learning on a budget: Opposite  
350 strategies suit high and low budgets. *arXiv preprint arXiv:2202.02794*, 2022.
- 351 [9] Qiang Hu, Yuejun Guo, Maxime Cordy, Xiaofei Xie, Wei Ma, Mike Papadakis, and Yves  
352 Le Traon. Towards exploring the limitations of active learning: An empirical study. In *2021*  
353 *36th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages  
354 917–929. IEEE, 2021.
- 355 [10] Yilin Ji, Daniel Kaestner, Oliver Wirth, and Christian Wressnegger. Randomness is the root  
356 of all evil: More reliable evaluation of deep active learning. In *Proceedings of the IEEE/CVF*  
357 *Winter Conference on Applications of Computer Vision*, pages 3943–3952, 2023.
- 358 [11] Jinhan Kim, Robert Feldt, and Shin Yoo. Guiding deep learning system testing using surprise  
359 adequacy. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*,  
360 pages 1039–1049. IEEE, 2019.
- 361 [12] Andreas Kirsch, Joost Van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch  
362 acquisition for deep bayesian active learning. *Advances in neural information processing sys-*  
363 *tems*, 32, 2019.
- 364 [13] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
365 2009.
- 366 [14] CURE Lab. Deep active learning with pytorch. [https://github.com/cure-lab/](https://github.com/cure-lab/deep-active-learning)  
367 *deep-active-learning*, 2022.
- 368 [15] Yu Li, Muxi Chen, Yannan Liu, Daojing He, and Qiang Xu. An empirical study on the efficacy  
369 of deep active learning for image classification. *arXiv preprint arXiv:2212.03088*, 2022.

- 370 [16] David Lowell, Zachary C Lipton, and Byron C Wallace. Practical obstacles to deploying active  
371 learning. *arXiv preprint arXiv:1807.04801*, 2018.
- 372 [17] Carsten Lüth, Till Bungert, Lukas Klein, and Paul Jaeger. Navigating the pitfalls of active  
373 learning evaluation: A systematic framework for meaningful performance assessment. *Ad-  
374 vances in Neural Information Processing Systems*, 36, 2024.
- 375 [18] Rishabh Misra. News category dataset. *arXiv preprint arXiv:2209.11429*, 2022.
- 376 [19] Prateek Munjal, Nasir Hayat, Munawar Hayat, Jamshid Sourati, and Shadab Khan. Towards ro-  
377 bust and reproducible active learning using neural networks. In *Proceedings of the IEEE/CVF  
378 Conference on Computer Vision and Pattern Recognition*, pages 223–232, 2022.
- 379 [20] Information Engineering Graduate Institute of Taiwan University. Libsvmtools.
- 380 [21] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for  
381 word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages  
382 1532–1543, 2014.
- 383 [22] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set  
384 approach. *arXiv preprint arXiv:1708.00489*, 2017.
- 385 [23] Dan Wang and Yi Shang. A new active labeling method for deep learning. In *2014 Interna-  
386 tional joint conference on neural networks (IJCNN)*, pages 112–119. IEEE, 2014.
- 387 [24] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for bench-  
388 marking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 389 [25] Donggeun Yoo and In So Kweon. Learning loss for active learning. In *Proceedings of the  
390 IEEE/CVF conference on computer vision and pattern recognition*, pages 93–102, 2019.
- 391 [26] Xueying Zhan, Huan Liu, Qing Li, and Antoni B Chan. A comparative survey: Benchmarking  
392 for pool-based active learning. In *IJCAI*, pages 4679–4686, 2021.
- 393 [27] Yilun Zhou, Adithya Renduchintala, Xian Li, Sida Wang, Yashar Mehdad, and Asish Ghoshal.  
394 Towards understanding the behaviors of optimal deep active learning algorithms. In *Interna-  
395 tional Conference on Artificial Intelligence and Statistics*, pages 1486–1494. PMLR, 2021.

## 396 A Acquisition Functions

397 **Uncertainty Sampling** tries to find the sample that the classifier is most uncertain about by  
398 computing heuristics of the class probabilities. For our benchmark, we use entropy and margin  
399 (a.k.a. best-vs-second-best) sampling.

400 **BALD** [12] applies the query-by-committee strategy of model ensembles to a single model by  
401 interpreting the classifier’s parameters as distributions and then sample multiple outputs from them  
402 via Monte-Carlo dropout.

403 **BADGE** [1] uses gradient embeddings of unlabeled points to select samples where the classifier  
404 is expected to change a lot. The higher the magnitude of the gradient the higher the expected  
405 improvement in model performance.

406 **Coreset** [22] employs K-Means clustering trying to cover the whole data distribution. Selects  
407 the unlabeled sample that is the furthest away from all cluster centers. Clustering is done in a  
408 semantically meaningful space by encoding the data with the current classifier  $\hat{y}$ . In this work, we  
409 use the greedy variant of Coreset.

410 **TypiClust** [8] relies on clustering similar to Coreset, but proposes a new measure called “Typical-  
411 ity” to select unlabeled samples. It selects points that are in the densest regions of clusters that do  
412 not contain labeled samples yet. Clustering is done in a semantically meaningful space by encoding  
413 the data with the current classifier  $\hat{y}$ . It has to be pointed out that TypiClust was designed for  
414 low-budget scenarios, but we think it is still worthwhile to test and compare this algorithm with  
415 higher budgets.

416 **Core-GCN** [3] TODO

417 **DSA/LSA** [11] TODO

### 418 Excluded Algorithms

419 **Learning Loss for AL** [25] Introduces an updated training of the classification model with an  
420 auxiliary loss and therefore cannot be compared fairly against classification models without this  
421 boosted training regime.

### 422 Reinforcement Learning Algorithms

423

## 424 B Difference of Ranks with 3 Repetitions

425 Table 4 and Table 5 follow the exact same computation of ranks that created the main result (Table  
426 3) with the only difference being a reduced number of runs per acquisition function. For each table  
427 we uniformly sampled 3 runs from the available 50 per acquisition function.

428 We can observe significant differences between the two tables:

429 **Purple:** A multitude of rank differences of acquisition functions for specific datasets, some as high  
430 as 4.7 ranks for TypiClust on the Splice dataset

431 **Olive:** Well separated acquisition functions in Tab. 5 (Margin and BADGE) are almost indistin-  
432 guishable in Tab 4

433 **Red:** BALD lost 2 places in the overall ranking and Entropy gained 2

434 Even though the overall ordering of acquisition functions stayed relatively unchanged due to the  
435 averaging across many datasets, each individual dataset was subject to drastic permutations. This  
436 highlights the need for many repetitions in AL experiments.

Table 4: Ranks of all acquisition functions per dataset. First random draw of 3 runs from the overall pool of 50.

	Splice	DNA	USPS	Cifar10	FMnist	TopV2	News	Unencoded	Encoded
Oracle	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	2.1
Margin	6.0	7.3	2.0	6.7	5.3	2.3	3.3	4.7	4.4
Badge	6.0	7.3	3.0	6.7	5.0	3.3	4.0	5.0	5.3
BALD	3.3	4.7	5.3	12.0	7.0	6.3	4.3	6.1	7.9
CoreGCN	8.7	3.7	10.7	6.3	5.3	4.0	7.7	6.6	9.1
DSA	8.3	6.3	7.7	7.7	4.3	6.7	6.7	6.8	6.1
LeastConf	10.0	12.0	8.0	3.0	4.3	9.3	2.3	7.0	6.7
LSA	5.7	6.7	5.3	6.7	10.7	7.7	7.0	7.1	6.3
Entropy	11.0	3.3	7.3	4.0	6.7	8.3	9.7	7.2	7.0
Random	7.7	8.7	5.3	8.0	11.0	8.0	9.0	8.2	6.3
Coreset	4.7	10.3	10.3	7.7	6.0	9.0	11.0	8.4	7.2
TypiClust	5.7	6.7	12.0	8.3	11.3	12.0	12.0	9.7	9.7

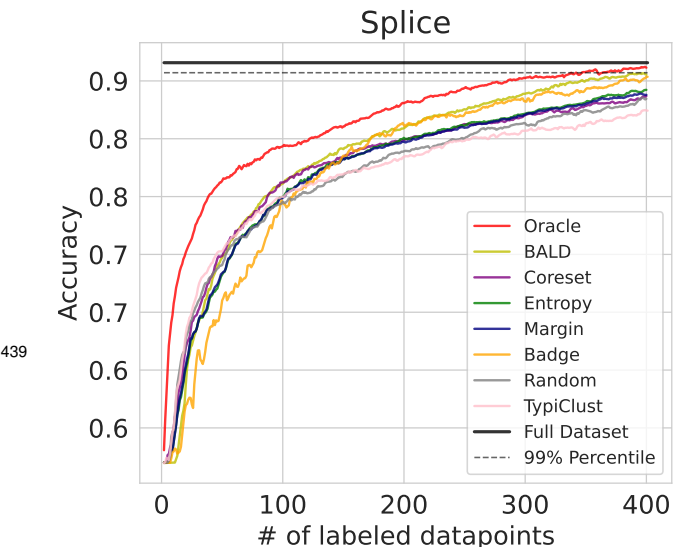
Table 5: Ranks of all acquisition functions per dataset. Second random draw of 3 runs from the overall pool of 50.

	Splice	DNA	USPS	Cifar10	FMnist	TopV2	News	Unencoded	Encoded
Oracle	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	2.4
Margin	6.0	3.3	2.0	5.7	2.0	2.0	4.3	3.6	3.8
Badge	6.0	9.0	3.0	3.0	5.7	3.7	3.3	4.8	4.9
CoreGCN	4.3	6.3	10.3	7.3	5.3	5.7	5.3	6.4	8.1
DSA	8.7	7.3	7.3	6.0	4.3	5.3	6.0	6.4	6.5
BALD	4.7	4.0	4.7	12.0	7.3	6.7	6.7	6.6	7.5
Entropy	6.7	4.7	7.7	5.3	5.0	7.3	9.3	6.6	6.8
LeastConf	7.7	10.0	8.3	3.3	6.0	8.7	3.0	6.7	7.3
LSA	7.7	5.3	6.0	9.0	11.0	9.0	7.3	7.9	7.5
Random	9.3	8.0	5.0	8.7	11.7	8.3	8.7	8.5	7.6
Coreset	6.0	10.7	10.7	8.0	8.3	8.3	11.0	9.0	6.3
TypiClust	10.0	8.3	12.0	8.7	10.3	12.0	12.0	10.5	9.4

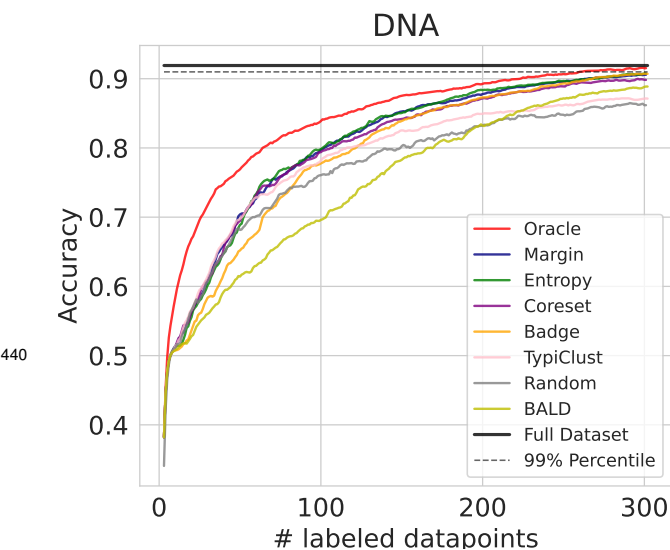
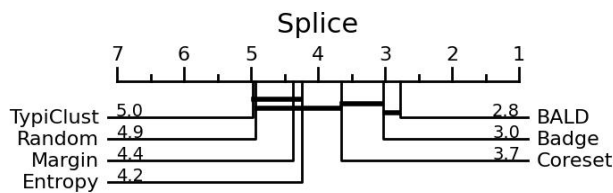


437 **C Ranks by Query Size**

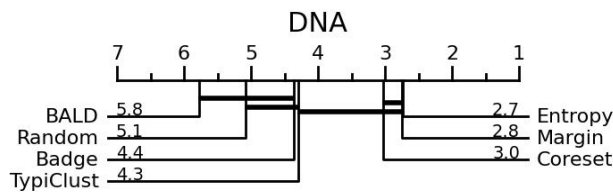
438 **D Individual Results**



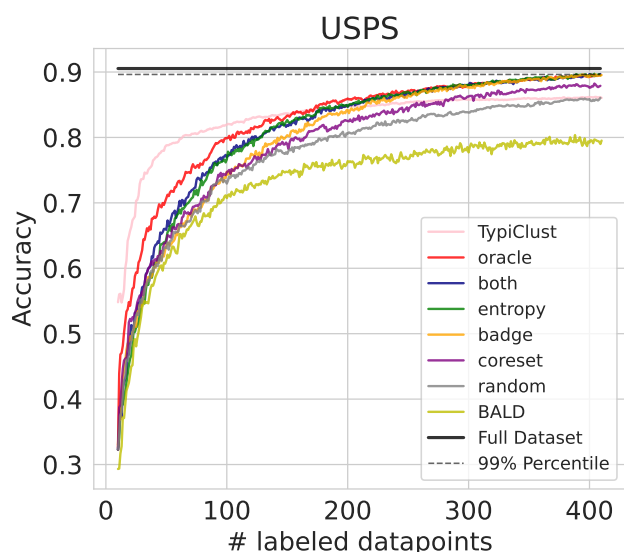
	Splice
Oracle	$0.811 \pm 0.010$
BALD	$0.785 \pm 0.013$
Coreset	$0.778 \pm 0.014$
Entropy	$0.774 \pm 0.016$
Margin	$0.773 \pm 0.016$
Badge	$0.770 \pm 0.016$
Random	$0.768 \pm 0.014$
TypiClust	$0.766 \pm 0.014$



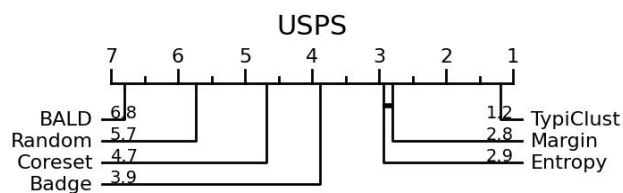
	DNA
Oracle	$0.842 \pm 0.021$
Margin	$0.807 \pm 0.035$
Entropy	$0.805 \pm 0.038$
Coreset	$0.796 \pm 0.028$
Badge	$0.789 \pm 0.056$
TypiClust	$0.788 \pm 0.036$
Random	$0.768 \pm 0.024$
BALD	$0.749 \pm 0.044$



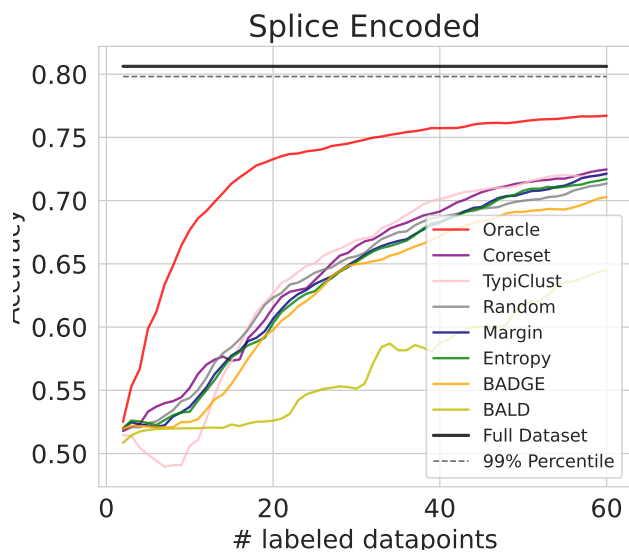
441



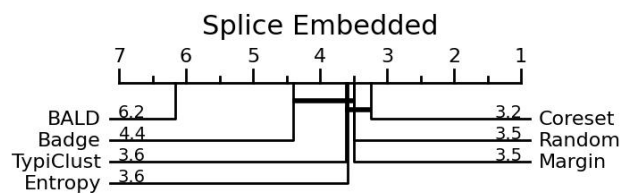
	USPS
TypiClust	$0.830 \pm 0.007$
Oracle	$0.823 \pm 0.011$
Margin	$0.809 \pm 0.013$
Entropy	$0.807 \pm 0.013$
Badge	$0.795 \pm 0.018$
Coreset	$0.788 \pm 0.017$
Random	$0.774 \pm 0.012$
BALD	$0.725 \pm 0.050$



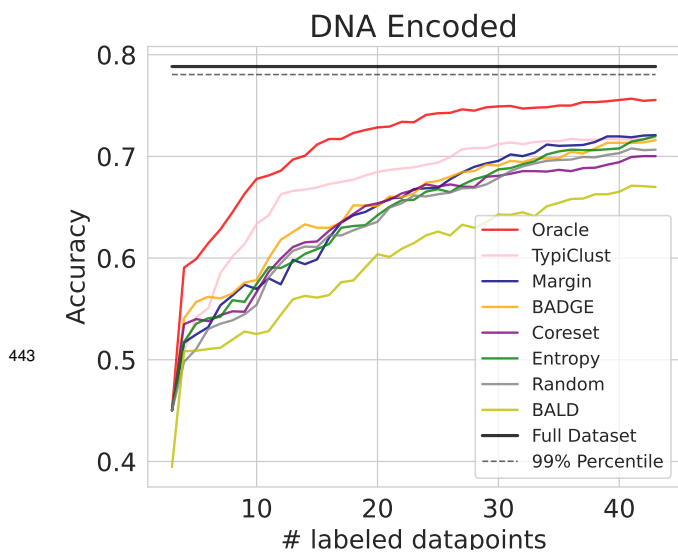
442



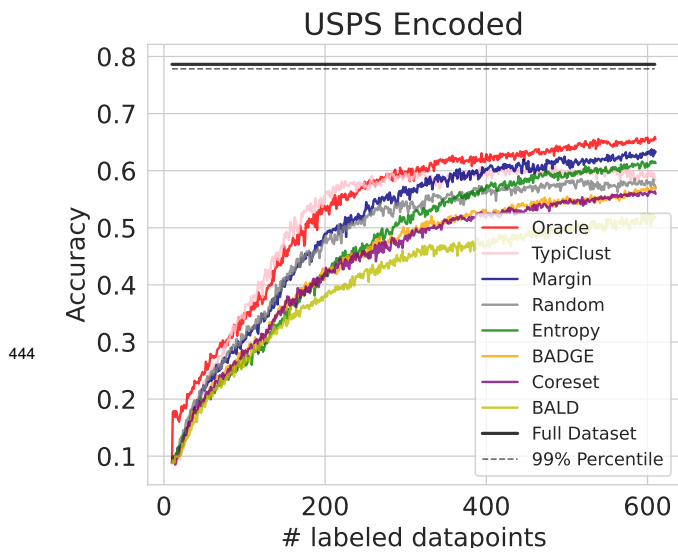
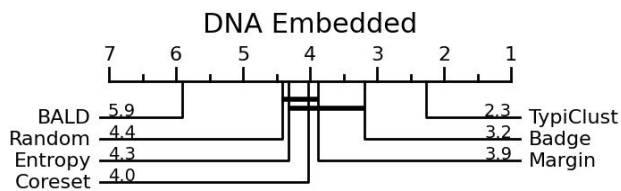
	SpliceEncoded
Oracle	$0.728 \pm 0.022$
Coreset	$0.648 \pm 0.027$
TypiClust	$0.645 \pm 0.042$
Random	$0.643 \pm 0.036$
Entropy	$0.636 \pm 0.033$
Margin	$0.636 \pm 0.033$
Badge	$0.627 \pm 0.040$
BALD	$0.565 \pm 0.049$



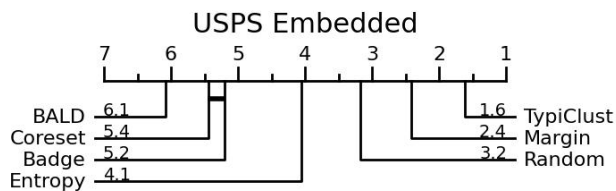


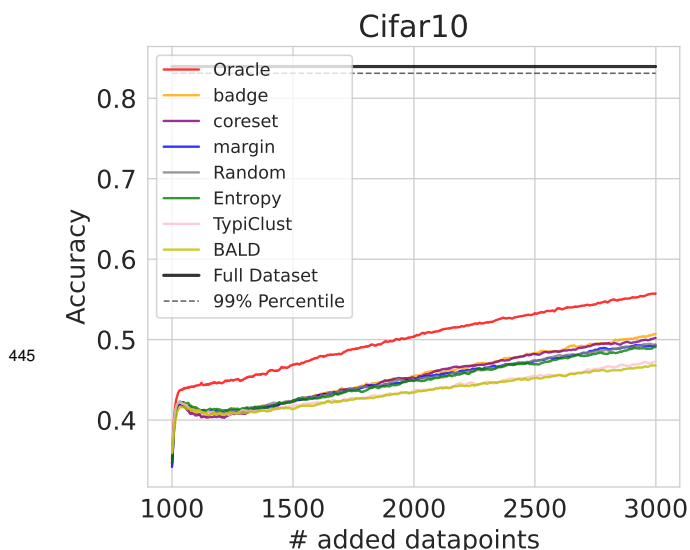


	DNAEncoded
Oracle	$0.709 \pm 0.023$
TypiClust	$0.672 \pm 0.029$
Margin	$0.648 \pm 0.047$
Badge	$0.647 \pm 0.037$
Coreset	$0.640 \pm 0.041$
Entropy	$0.629 \pm 0.062$
Random	$0.626 \pm 0.035$
BALD	$0.594 \pm 0.039$

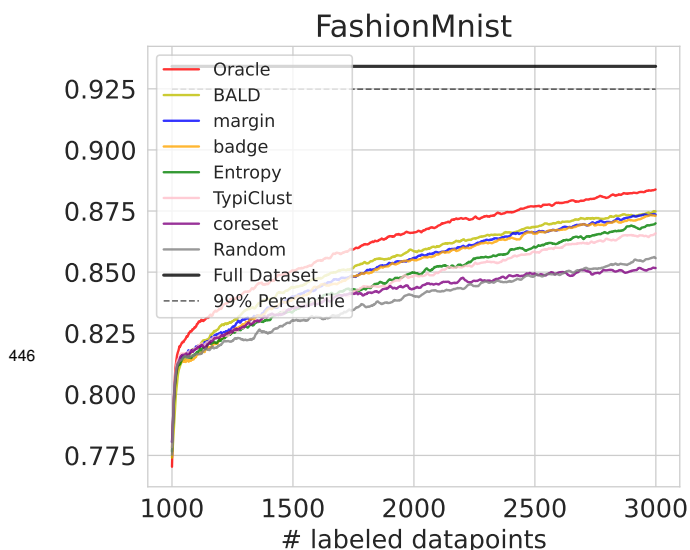
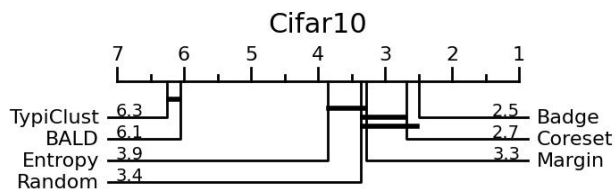


	USPSEncoded
Oracle	$0.522 \pm 0.021$
TypiClust	$0.507 \pm 0.025$
Margin	$0.496 \pm 0.030$
Random	$0.468 \pm 0.025$
Entropy	$0.459 \pm 0.021$
Badge	$0.440 \pm 0.026$
Coreset	$0.435 \pm 0.027$
BALD	$0.402 \pm 0.052$

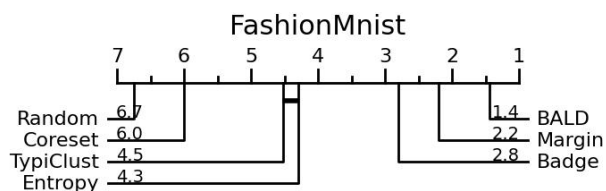




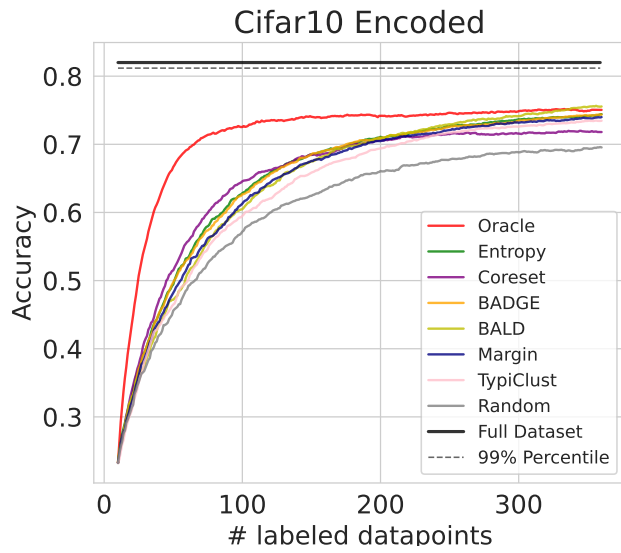
	Cifar10
Oracle	$0.500 \pm 0.010$
Badge	$0.453 \pm 0.012$
Coreset	$0.453 \pm 0.009$
Margin	$0.451 \pm 0.010$
Random	$0.450 \pm 0.012$
Entropy	$0.449 \pm 0.010$
TypiClust	$0.436 \pm 0.010$
BALD	$0.436 \pm 0.010$



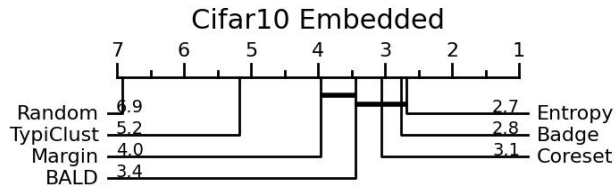
	FashionMnist
Oracle	$0.862 \pm 0.003$
BALD	$0.854 \pm 0.003$
Margin	$0.851 \pm 0.003$
Badge	$0.851 \pm 0.003$
Entropy	$0.847 \pm 0.004$
TypiClust	$0.846 \pm 0.004$
Coreset	$0.840 \pm 0.004$
Random	$0.837 \pm 0.004$



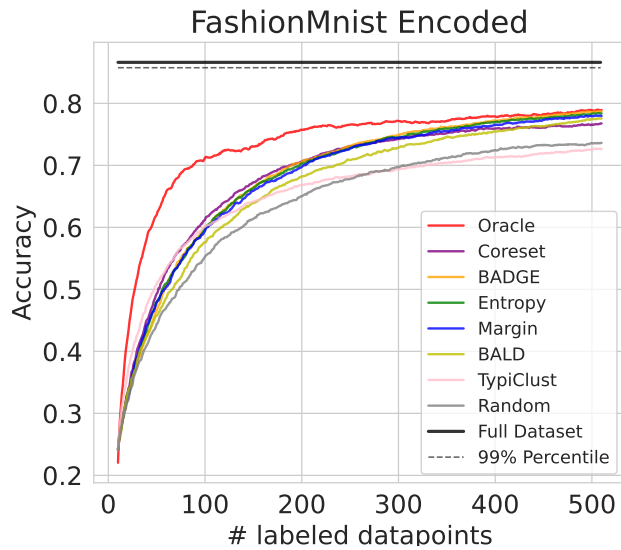
447



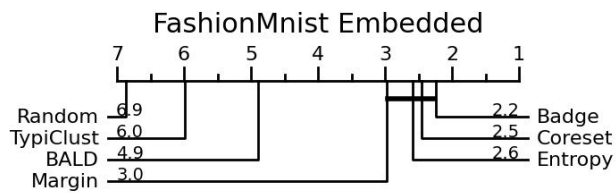
	Cifar10Encoded
Oracle	$0.714 \pm 0.007$
Entropy	$0.654 \pm 0.013$
Coreset	$0.653 \pm 0.012$
Badge	$0.653 \pm 0.012$
BALD	$0.650 \pm 0.016$
Margin	$0.647 \pm 0.012$
TypiClust	$0.636 \pm 0.009$
Random	$0.607 \pm 0.013$

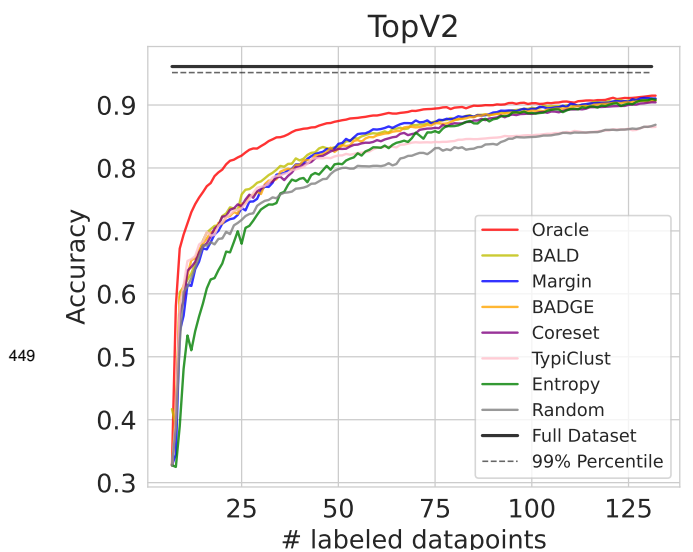


448

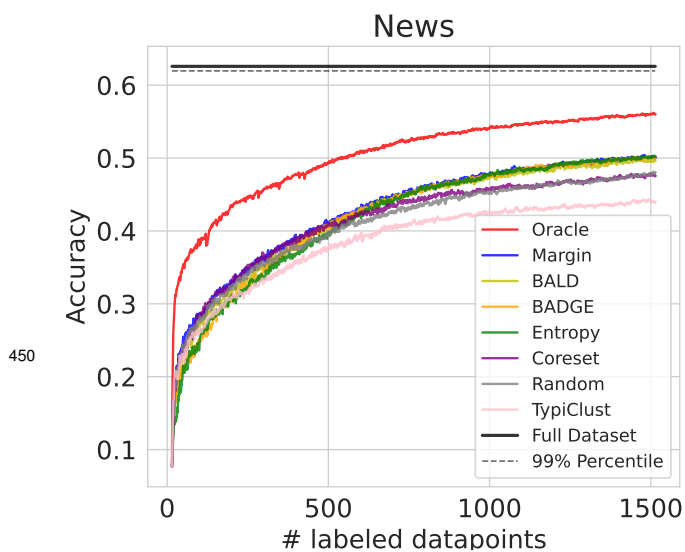
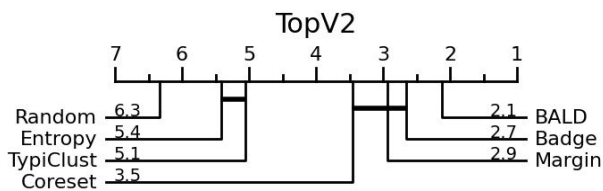


	FashionMnistEncoded
Oracle	$0.732 \pm 0.006$
Coreset	$0.686 \pm 0.008$
Badge	$0.685 \pm 0.008$
Entropy	$0.684 \pm 0.009$
Margin	$0.682 \pm 0.011$
BALD	$0.668 \pm 0.009$
TypiClust	$0.652 \pm 0.009$
Random	$0.640 \pm 0.011$

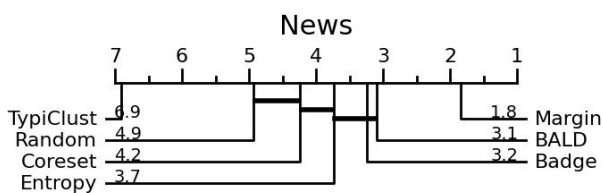


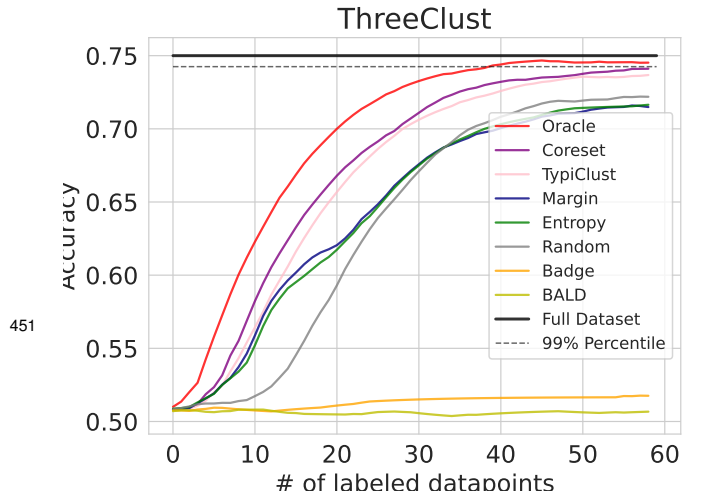


	TopV2
Oracle	$0.862 \pm 0.006$
BALD	$0.831 \pm 0.013$
Badge	$0.826 \pm 0.015$
Coreset	$0.823 \pm 0.016$
Margin	$0.822 \pm 0.015$
TypiClust	$0.805 \pm 0.015$
Entropy	$0.801 \pm 0.025$
Random	$0.787 \pm 0.015$

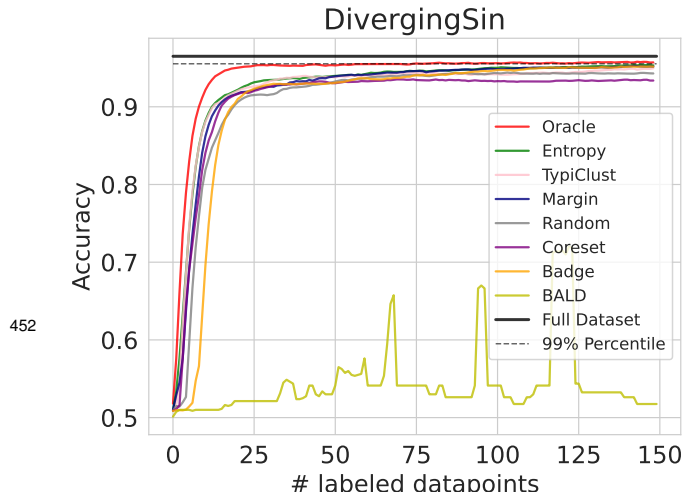
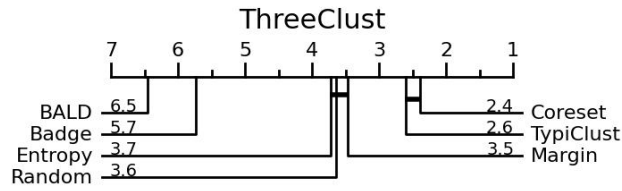


	News
Oracle	$0.502 \pm 0.005$
Margin	$0.427 \pm 0.007$
BALD	$0.421 \pm 0.008$
Badge	$0.420 \pm 0.011$
Entropy	$0.416 \pm 0.010$
Coreset	$0.415 \pm 0.011$
Random	$0.409 \pm 0.008$
TypiClust	$0.385 \pm 0.010$

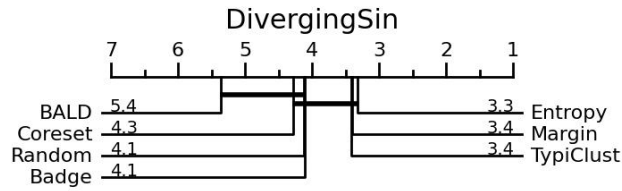




	ThreeClust
Oracle	$0.722 \pm 0.097$
Coreset	$0.698 \pm 0.058$
TypiClust	$0.697 \pm 0.055$
Entropy	$0.682 \pm 0.098$
Random	$0.672 \pm 0.067$
Margin	$0.669 \pm 0.095$
Badge	$0.524 \pm 0.086$
BALD	$0.507 \pm 0.050$



	DivergingSin
Oracle	$0.948 \pm 0.198$
Entropy	$0.936 \pm 0.202$
TypiClust	$0.930 \pm 0.196$
Margin	$0.929 \pm 0.201$
Random	$0.919 \pm 0.191$
Badge	$0.914 \pm 0.202$
Coreset	$0.914 \pm 0.197$
BALD	$0.661 \pm 0.167$



## 453 E Seeding Strategy

454 We aim to provide an experimental setup that is fully reproducible independent of the dataset, classi-  
 455 fication model, or AL algorithm used. For a fair comparison of two AL algorithms, both algorithms  
 456 need to receive equal starting conditions in terms of train/validation split, initialization of classifier,

and even the state of minor systems like the optimizer or mini-batch sampler. Even though different implementations might have their own solution to some of these problems, only [10] has described and implemented a fully reproducible pipeline for AL evaluation. The term reproducibility in this work is used as a synonym not only for the reproducibility of an experiment (a final result given a seed), but also the reproducibility of all subsystems independent of each other. The seed for one subsystem should always reproduce the behavior of this subsystem independent of all other subsystems and their seeds. The main obstacle for ensuring reproducibility is the seeding utility in PyTorch, Tensorflow, and other frameworks, whose default choice is a single global seed. Since many subsystems draw random numbers from this seed, all of them influence each other to a point where a single additional draw can completely change the model initialization, data split or the order of training batches. Even though some workarounds exist, e.g. re-setting the seed multiple times, this problem is not limited to the initialization phase, but also extends to the AL iterations and the systems within. We propose an implementation that creates separate Random Number Generators (RNGs) for each of these systems to ensure equal testing conditions even when the AL algorithm, dataset, or classifier changes. We hypothesize that the insufficient setup with global seeds contributes to the ongoing problem of inconsistent results of AL algorithms in different papers.

In summary, we introduce three different seeds:  $s_\Omega$  for the AL algorithm,  $s_\mathcal{D}$  for dataset splitting and mini-batch sampling, and  $s_\theta$  for model initialization and sampling of dropout masks. Unless stated otherwise, we will keep  $s_\Omega$  fixed, while  $s_\mathcal{D}$  and  $s_\theta$  are incremented by 1 between restarts to introduce stochasticity into our framework. Some algorithms require a subsample to be drawn from  $\mathcal{U}$  in order to reduce the computational cost in each iteration, while others need access to the full unlabeled pool (e.g. for effective clustering). If a subsample is required, it will be drawn from  $s_\Omega$  and therefore will not influence other systems in the experiments. For each algorithm, we decided if subsampling is required based on our available hardware, but decided against setting a fixed time limit per experiment, since this would introduce unnecessary complexity into the benchmark. An overview of selected hyperparameters per AL algorithm can be found in Appendix H.

**Note:** Even though we decoupled the subsystems via the described seeds, the subsystems can still influence each other in a practical sense. For example, keeping  $s_\mathcal{D}$  fixed does not mean that always the same sequence of samples from  $\mathcal{U}$  (if subsamples are drawn) are shown to all acquisition functions. This is practically impossible, as different acquisition functions pick different  $x^{(i)}$ . However, the hypothetical **tree** of all possible sequences of samples from  $\mathcal{U}$  remains the same, granting every acquisition function equal possibilities.

## F Number of Restarts Ablation for "Training from Scratch"

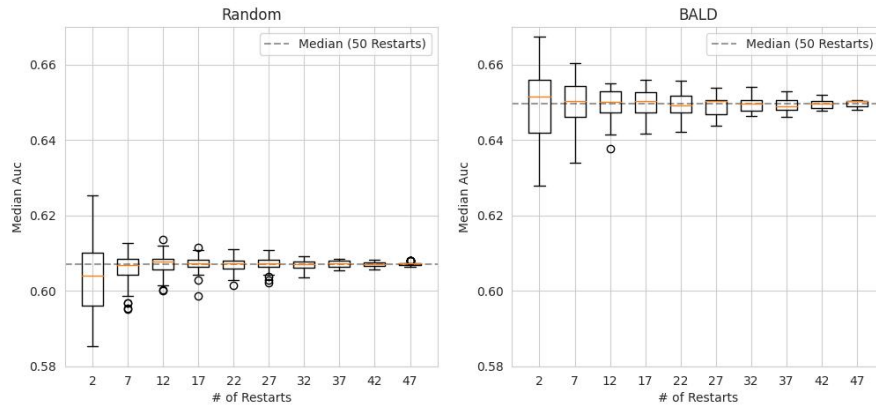


Figure 5: Created by the protocol described in Sec. 4.1 for the Encoded Cifar10 dataset, with the classifier being trained from scratch. We observe similar characteristics, where a low number of restarts leads to extreme variance in the reported results.

## G Oracle Curve Forecasting

TODO

## H Hyperparameters per AL Algorithm

Table 6: Selected hyperparameters for all tested acquisition functions. Last column indicates the source of our implementation.

Algorithm	Sample Size	Other	Source
BADGE	100	Dropout Trials: 5	Based on [1, 14]
BALD	100		Based on [4]
Coreset	8000		Own
TypiClust	10000	Min Cluster Size: 5 Max # Clusters: 500	Based on [8]
Margin	8000		Own
Entropy	8000		Own

## I Hyperparameters and Preprocessing per Dataset

For all our datasets we use the pre-defined train/test splits, if given. In the remaining cases, we define test sets upfront and store them into separate files to keep them fixed across all experiments. The validation set is split in the experiment run itself and depends on the dataset-seed.

**Tabular:** We use **Splice**, **DNA** and **USPS** from LibSVMTools [20]. All three datasets are normalized between [0, 1].

**Image:** We use **FashionMNIST** [24] and **Cifar10** [13], since both are widely used in AL literature. Both datasets are normalized according to their standard protocols.

**Text:** We use **News Category** [18] and **TopV2** [6]. For News Category we use the 15 most common categories as indicated by its Kaggle site. We additionally drop sentences above 80 words to reduce the padding needed (retaining 99,86% of the data). For TopV2, we are only using the "alarm" domain. Both datasets are encoded with pre-trained GloVe (Common Crawl 840B Tokens) embeddings [21]. Since neither dataset provided a fixed test set, we randomly split 7000 datapoints into a test set.

Dataset	Seed Set	Budget	Val Split
Splice	1	400	0.2
SpliceEnc.	1	60	0.2
DNA	1	300	0.2
DNAEnc	1	40	0.2
USPS	1	400	0.2
USPSEnc	1	600	0.2
FashionMnist	100	2000	0.04
FashionMnistEnc	1	500	0.04
Cifar10	100	2000	0.04
Cifar10Enc	1	350	0.04
TopV2	1	125	0.25
News	1	1500	0.03

Table 7: Size of the seed set is given by number of labeled sample per class.

Dataset	Classifier	Optimizer	LR	Weight Decay	Dropout	Batch Size
Splice	[24, 12]	NAdam	1.2e-3	5.9e-5	0	43
SpliceEnc.	linear	NAdam	6.2e-4	5.9e-6	0	64
DNA	[24, 12]	NAdam	3.9e-2	3.6e-5	0	64
DNAEnc	linear	NAdam	1.6e-3	4e-4	0	64
USPS	[24, 12]	Adam	8.1e-3	1.5e-6	0	43
USPSEnc	linear	NAdam	7.8e-3	1.9e-6	0	64
FashionMnist	ResNet18	NAdam	1e-3	0	0	64
FashionMnistEnc	linear	Adam	1.6e-3	1e-5	5e-2	64
Cifar10	ResNet18	NAdam	1e-3	0	0	64
Cifar10Enc	linear	NAdam	1.7e-3	2.3e-5	0	64
TopV2	BiLSTM	NAdam	1.5e-3	1.7e-7	5e-2	64
News	BiLSTM	NAdam	1.5e-3	1.7e-7	5e-2	64

Table 8: Classifier architectures and optimized hyperparameters per dataset. Numbers in brackets signify a MLP with corresponding hidden layers.

## 507 J AL Pseudocode

---

### Algorithm 1 Active Learning Loop

---

**Require:**  $\mathcal{L}, \mathcal{U}, \mathcal{D}_{\text{test}}, \text{Train}, \text{Seed}, \hat{y}$

**Require:**  $\Omega$

```

1:  $\mathcal{L}^{(1)} \leftarrow \text{Seed}(\mathcal{U})$ 
2:  $\mathcal{U}^{(1)} \leftarrow \mathcal{U}$ 
3: for  $i := 1 \dots B$  do
4:    $\text{acc}^{(i)} \leftarrow \text{Train}(\mathcal{L}^{(i)})$ 
5:    $a^{(i)} \leftarrow \Omega(\mathcal{U}^{(i)})$ 
6:    $\mathcal{L}^{(i+1)} \leftarrow \mathcal{L}^{(i)} \cup \{(\mathcal{U}_a^{(i)}, A(\mathcal{U}_a^{(i)}))\}$ 
7:    $\mathcal{U}^{(i+1)} \leftarrow \mathcal{U}^{(i)} \setminus \{\mathcal{U}_a^{(i)}\}$ 
8: return  $\frac{1}{B} \sum_{i=1}^B \text{acc}^{(i)}$ 

```

---

▷ Acquisition Function

▷ Create the initial labeled set

---

### Algorithm 2 Retrain

---

**Require:**  $\mathcal{L}, \mathcal{D}_{\text{val}}, \mathcal{D}_{\text{test}}$

**Require:**  $\hat{y}, e_{\text{max}}$

```

1:  $\text{loss}^* \leftarrow \infty$ 
2: for  $i := 1 \dots e_{\text{max}}$  do
3:    $\hat{y}_{i+1} \leftarrow \hat{y}_i - \eta \nabla_{\hat{y}} \ell(\mathcal{L}, \hat{y})$ 
4:    $\text{loss}_i \leftarrow \ell(\mathcal{D}^{\text{val}}, \hat{y})$ 
5:   if  $\text{loss}_i < \text{loss}^*$  then
6:      $\text{loss}^* \leftarrow \text{loss}_i$ 
7:   else
8:     Break
9: return  $\text{Acc}(\mathcal{D}^{\text{test}}, \hat{y})$ 

```

---



---

**Algorithm 3** Acquire Oracle  $\Omega$ 

---

**Require:**  $\mathcal{U}, \mathcal{L}, A, \mathcal{D}_{\text{test}}, \tau, \hat{y}_\theta$ **Require:** Train, Margin, Acc

```
1:  $\text{acc}^0 \leftarrow \text{acc}^* \leftarrow \text{Acc}(\mathcal{D}_{\text{test}}, \hat{y}_\theta)$ 
2: for  $k := 1 \dots \tau$  do
3:    $u_k = \text{unif}(\mathcal{U})$ 
4:    $\mathcal{L}' \leftarrow \mathcal{L}^{(i)} \cup \{(u_k, A(u_k))\}$ 
5:    $\hat{y}'_\theta \leftarrow \text{Train}(\mathcal{L}', \hat{y}_\theta)$ 
6:    $\text{acc}' \leftarrow \text{Acc}(\mathcal{D}_{\text{test}}, \hat{y}'_\theta)$ 
7:   if  $\text{acc}' > \text{acc}^*$  then
8:      $\text{acc}^* \leftarrow \text{acc}'$ 
9:      $u^* \leftarrow u_k$ 
10: if  $\text{acc}^0 = \text{acc}^*$  then
11:    $u^* \leftarrow \text{Margin}(\mathcal{U}, \hat{y}_\theta)$ 
return  $u^*$ 
```

---

508 Alg. 3 replaces the acquisition function  $\Omega$  in the AL loop (Alg. J line 5).