# **Towards Robust and Reproducible Active Learning Using Neural Networks**

Prateek Munjal<sup>1</sup> Nasir Hayat<sup>1</sup> Munawar Hayat<sup>1</sup> Jamshid Sourati<sup>2</sup> Shadab Khan<sup>1</sup>

## **Abstract**

Active learning (AL) is a promising ML paradigm that has the potential to parse through large unlabeled data and help reduce annotation cost in domains where labeling entire data can be prohibitive. Recently proposed neural network based AL methods use different heuristics to accomplish this goal. In this study, we show that recent AL methods offer a gain over random baseline under a brittle combination of experimental conditions. We demonstrate that such marginal gains vanish when experimental factors are changed, leading to reproducibility issues and suggesting that AL methods lack robustness. We also observe that with a properly tuned model, which employs recently proposed regularization techniques, the performance significantly improves for all AL methods including the random sampling baseline, and performance differences among the AL methods become negligible. Based on these observations, we suggest a set of experiments that are critical to assess the true effectiveness of an AL method. To facilitate these experiments we also present an open source toolkit. We believe our findings and recommendations will help advance reproducible research in robust AL using neural networks.

*Key Abbreviations*: AL: Active Learning, RA: Random Augmentation, SWA: Stochastic Weight Averaging, SS: Shake-Shake, RSB: Random Sampling Baseline

#### 1. Introduction

Active learning (AL) is a machine learning paradigm that promises to help reduce the burden of data annotation by intelligently selecting a subset of informative samples from a large pool of unlabeled data that are relatively more conducive for learning. In AL, a model trained with a small amount of labeled seed data is used to parse through the

unlabeled data to select the subset that should be sent to an annotator (called oracle in AL literature). To select such a subset, AL methods rely on exploiting the latent-space structure of samples, model uncertainty, or other such heuristics. The promise of reducing annotation cost has brought a surge in recent AL research (Sinha et al., 2019), (Sener & Savarese, 2018), (Beluch et al., 2018), (Gal et al., 2017), (Kirsch et al., 2019), (Tran et al., 2019), (Yoo & Kweon, 2019), and with it, a few outstanding issues.

First, the results reported for RSB vary significantly between studies. For example, using 20% labeled data of CI-FAR10, the difference between RSB performance reported by (Yoo & Kweon, 2019) and (Tran et al., 2019) is 13\% under identical settings. **Second**, the results reported for the same AL method can vary across studies: using VGG16 (Simonyan & Zisserman, 2014) on CIFAR100 (Krizhevsky & Hinton, 2009) with 40% labeled data, (Sener & Savarese, 2018) reports  $\approx 55\%$  classification accuracy whereas (Sinha et al., 2019) reports 47.01% for (Sener & Savarese, 2018). **Third**, recent AL studies have been inconsistent with each other. For example, (Sener & Savarese, 2018) and (Ducoffe & Precioso, 2018) state that diversity-based AL methods consistently outperform uncertainty-based methods, which were found to be worse than the random sampling baseline (RSB). In contrast, recent developments in uncertainty based studies (Yoo & Kweon, 2019) suggest otherwise.

In addition to these issues, results using a new AL method are often reported on simplistic datasets and tested under limited experimental conditions, with an underlying assumption that the relative performance gains using an AL method would be maintained under changes in the experimental conditions. These issues with reporting of AL results has spurred a recent interest in benchmarking of AL methods and recent NLP and computer vision studies have raised a number of interesting questions (Lowell et al., 2018), (Prabhu et al., 2019), (Mittal et al., 2019). With the goal of improving the reproducibility and robustness of AL methods, in this study we evaluate the performance of these methods for image classification compared to a RSB in a fair experimental environment. The contributions of this study are as follows.

Contributions: Through a comprehensive set of experiments performed using our PyTorch-based AL evalua-

<sup>&</sup>lt;sup>1</sup>Inception Institute of AI, Abu Dhabi, UAE <sup>2</sup>University of Chicago, Chicago, IL, USA. Correspondence to: Shadab Khan <skhan.shadab@gmail.com>.

tion toolkit1 we compare different AL methods including state-of-the-art diversity-based, uncertainty-based, and committee-based methods (Sinha et al., 2019), (Sener & Savarese, 2018), (Beluch et al., 2018), (Gal et al., 2017) and a well-tuned RSB. We demonstrate that: 1) results with our RSB are higher across a range of experiments than previously reported, 2) state-of-the-art AL methods achieve a marginal gain over our RSB under narrow combination of experimental conditions (e.g. a specific architecture), which vanishes with changes in experimental conditions (e.g. using a different architecture for classifier), 3) variance in evaluation metric (accuracy) across repeated runs on the same set of data, or on different fold of initial labeled data, can lead to incorrect conclusions where accuracy gain using an AL method may be observed within the margin of error of accuracy measurement, 4) a bit surprisingly, our experiments also show that these performance gains vanish when the neural networks are well-regularized, and none of the evaluated AL methods performs better than our RSB 5) the variance in accuracy achieved using AL methods is substantially lower in consistent repeated training runs with a well-regularized model, suggesting that such a training regime is unlikely to effect misleading results in AL experiments, 6) finally, we conclude the paper with a set of guidelines on experimental evaluation of a new AL method, and provide a PyTorch-based AL toolkit to facilitate this.

# 2. Pool-Based Active Learning Methods

Contemporary pool-based AL methods can be broadly classified into: (i) uncertainty based (Sinha et al., 2019), (Gal et al., 2017), (Kirsch et al., 2019), (ii) diversity based (Sener & Savarese, 2018), (Ducoffe & Precioso, 2018), and (iii) committee based (Beluch et al., 2018). AL methods also differ in other aspects, for example, some AL methods use the task model (e.g. model trained for image classification) within their sampling function (Gal et al., 2017), (Sener & Savarese, 2018), where as others use different models for task and sampling functions (Sinha et al., 2019), (Beluch et al., 2018). These methods are discussed in detail next.

**Notations**: Starting with an initial set of labeled data  $L_0^0 = \{(x_i, y_i)\}_{i=1}^{N_L}$  and a large pool of unlabeled data  $U_0^0 = \{x_i\}_{i=1}^{N_U}$ , pool-based AL methods train a model  $\Phi_0$ . A sampling function  $\Psi(L_0^0, U_0^0, \Phi_0)$  then evaluates  $x_i \in U_0$ , and selects k (budget size) **samples** to be labeled by an oracle. The selected samples with **oracle-annotated** labels are then added to  $L_0^0$ , resulting in an extended  $L_0^1$  labeled set, which is then used to **retrain**  $\Phi$ . This cycle of **sample-annotate-train** is repeated until the sampling budget is exhausted or a satisficing metric is achieved. AL sampling

functions evaluated in this study are outlined next.

#### 2.1. Model Uncertainty on Output (UC)

The method in (Lewis & Gale, 1994) ranks the unlabeled datapoints,  $x_i \in U$  in a descending order based on their scores given by  $\max_j \Phi(x_i); j \in \{1 \dots C\}$ , where C is the number of classes, and chose the top k samples. Typically this approach focuses on the samples in U for which the softmax classifier is least confident.

### 2.2. Deep Bayesian Active Learning (DBAL)

(Gal et al., 2017) train the model  $\Phi$  with dropout layers and use Monte carlo dropout to approximate the sampling from posterior. For our experiments, we used the two most reported acquisitions *i.e.*, max entropy and Bayesian Active Learning by Disagreement (BALD). The max entropy method selects the top k datapoints having maximum entropy ( $\arg\max_i \mathbb{H}[P(\mathbf{y}|x_i)]; \forall x_i \in U_0$ ) where the posterior is given by,  $P(\mathbf{y}|x_i) = \sum_{j=1}^T \frac{1}{T}P(\mathbf{y}|x_i,\phi_j)$ ; where T denotes number of forward passes through the model,  $\Phi$ . BALD selects the top k samples that increase the information gain over the model parameters i.e.,  $\arg\max_i \mathbb{I}[P(\mathbf{y},\Phi|x_i,L_0)]; \forall x_i \in U_0$ . We implement DBAL as described in (Gal et al., 2017) where probability terms in information gain is evaluated using above equation.

#### 2.3. Center of Gravity (CoG)

Uncertainty in unlabeled datapoints is estimated in terms of the euclidean distance from the centre of gravity  $(z_{\rm cog})$  in the latent space. We define the COG as:  $z_{\rm cog} = \sum_{i=1}^{N_L+N_U} \frac{\Phi^l(x_i)}{|N_L+N_U|}$ , where  $\Phi^l(x_i)$  denotes the  $l^{\rm th}$  layer activations of the model  $\Phi$  for  $x_i$ . Using this distance estimate, we select the top k farthest datapoints from CoG. For our experiments, we use the penultimate layer activations.

#### 2.4. Coreset

(Sener & Savarese, 2018) exploit the geometry of datapoints and choose samples that provide a cover to all datapoints. Essentially, their algorithm tries to find a set of points (cover-points), such that distance of any datapoint from its nearest cover-point is minimized. They proposed two sub-optimal but efficient solutions to this NP-Hard problem: coreset-greedy and coreset-MIP (Mixed Integer programming), coreset-greedy is used to initialize coreset-MIP. For our experiments, following (Yoo & Kweon, 2019), we implement coreset-greedy since it achieves comparable performance while being significantly compute efficient.

<sup>&</sup>lt;sup>1</sup>AL Toolkit will be released on GitHub. To get access to pre-release version, please contact Shadab Khan at skhan.shadab@gmail.com

#### 2.5. Variational Adversarial Active Learning (VAAL)

(Sinha et al., 2019) combined a VAE (Kingma & Welling, 2013) and a discriminator (Goodfellow et al., 2014) to learn a metric for AL sampling. VAE encoder is trained on both L and U, and the discriminator is trained on the latent space representations of L and U to distinguish between seen (L) and unseen (U) images. Sampling function selects samples from U with lowest discriminator confidence (to be seen) as measured by output of discriminator's softmax. Effectively, samples that are most likely to be unseen based on the discriminator's output are chosen.

## **Algorithm 1** AL Training Schedule

```
1: Input AL_{iter}, Budget size k and Oracle, A
  2: Split \mathcal{D} \to \{T_r, T_s, V\}
  3: Split T_r \to \{L_0^0, U_0^0\}
  4: Train a base classifier, \mathcal{B} using only L_0^0
  5: \phi = \mathcal{B}
 6: while i \in \{0 ... AL_{iter}\} do
7: sample \{x_j\}_{j=1}^k \in U_i \text{ using } \Psi(L_0^i, U_0^i, \phi)
8: \{x_j, y_j\}_{j=1}^k \leftarrow \{x_j, \mathcal{A}(x_j)\}_{j=1}^k
9: i \leftarrow i+1
           L_0^i \leftarrow L_0^i \cup \{x_j, y_j\}_{j=1}^k \\ U_0^i \leftarrow U_0^i \setminus \{x_j, y_j\}_{j=1}^k
10:
11:
            \phi \leftarrowInitialize randomly
12:
13:
            while convergence do
14:
                 Train \phi using only L_0^i
15:
            end while
16: end while
```

#### 2.6. Ensemble Variance Ratio Learning

Proposed by (Beluch et al., 2018), this is a query-bycommittee (QBC) method that uses a variance ratio computed by  $v = 1 - f_m/N$  to select the sample set with the largest dispersion (v), where N is the number of committee members (CNNs), and  $f_m$  is the number of predictions in the modal class category. Variance ratio lies in 0-1 range and can be treated as an uncertainty measure. We note that it is possible to formulate several AL strategies using the ensemble e.g. BALD, max-entropy, etc. Variance ratio was chosen for this study because it was shown by authors to lead to superior results. For training the CNN ensembles, we train 5 models with VGG16 architecture but a different random initialization. Further, following (Beluch et al., 2018), the ensembles are used only for sample set selection, a separate task classifier is trained in fully-supervised manner to do image classification.

## 3. Regularization and Active Learning

In a ML training pipeline comprising data–model–metric and training tricks, regularization can be introduced in sev-

eral forms. In neural networks, regularization is commonly applied using parameter norm penalty (metric), dropout (model), or using standard data augmentation techniques such as horizontal flips and random crops (data). However, parameter norm penalty coefficients are not easy to tune and dropout effectively reduces model capacity to reduce the extent of over-fitting on the training data, and requires the drop probability to be tuned. On the other hand, several recent studies in semi-supervised learning (SSL) have shown promising new ways of regularizing neural networks to achieve impressive gains. While it isn't surprising that these regularization techniques help reduce generalization error, most AL studies have overlooked them. We believe this is because of a reasonable assumption that if an AL method works better than random sampling, then its relative advantage should be maintained when newer regularization techniques and training tricks are used. Since regularization is critical for low-data training regime of AL where the massively-overparameterized model can easily overfit to the limited training data, we investigate the validity of such assumptions by applying regularization techniques to the entire data-model-metric chain of neural network training.

Specifically, we employ parameter norm penalty, random augmentation (RA) (Cubuk et al., 2019), stochastic weighted averaging (SWA) (Izmailov et al., 2018), and shake-shake (SS) (Gastaldi, 2017). In RA, a sequence of nrandomly chosen image transforms are sequentially applied to the training data, with a randomly chosen distortion magnitude (m) which picks a value between two extremes. For details of extreme values used for each augmentation choice, we refer the reader to work of (Cubuk et al., 2018). SWA is applied on the model by first saving e snapshots of model during the time-course of optimization, and then averaging the snapshots as a post-processing step. For SS experiments, we utilize the publicly available pytorch implementation<sup>2</sup>. The hyper-parameters associated with these techniques as well as experiments and results with regularization applied to neural network training with AL-selected sample sets are discussed in Sec. 5.3.

## 4. Implementation Details

We perform experiments on CIFAR10, CIFAR100, and ImageNet by following the training schedule summarized in Alg. 1. Given a dataset  $\mathcal{D}$ , we split it into train  $(T_r)$ , validation (V), and test  $(T_s)$  sets. The train set is further divided into the initial labeled  $(L_0)$  and unlabeled  $(U_0)$  sets. A base classifier  $\mathcal{B}$  is first trained, followed by iterations of sample-annotate-train process using various AL methods. Model selection is done by choosing the best performing model on the validation set. For a fair comparison, a consistent set of experimental settings is used across all methods.

https://github.com/hysts/pytorch\_shake\_shake

Dataset-specific training details are discussed next.

Learning rate (lr) and weight decay (wd) were tuned using grid search, and set as follows for individual datasets. CIFAR10: optimizer=Adam (Kingma & Ba, 2015), lr = $5e^{-4}$ ,  $wd = 5e^{-4}$ , input pre-processed using random horizontal flip (p = 0.5) and normalization (divide by 255). **CIFAR100**: optimizer=Adam,  $lr = 5e^{-4}$  and wd = 0 for AL iterations and  $lr = 5e^{-5}$  and wd = 0 for base classifier that was trained on  $L_0$ , input pre-processed using random crop (pad=4) followed by horizontal flip (p = 0.5) and normalization (divide by 255). ImageNet: optimizer=SGD,  $wd = 3e^{-4}$ . We train the base classifier on  $L_0$  for 200 epochs where lr = 0.1 with a linear warm-up schedule (for first 5 epochs) followed by decaying the lr by a factor of 10 on epoch number: {140, 160, 180}. For AL iterations we fine-tune the best model (picked by validation set accuracy) from previous iteration for 100 epochs where  $lr = 1e^{-2}$ which gets decayed by a factor of 10 on epoch number:  $\{35, 55, 80\}$ . Further, we choose the best model based on a realistically small validation set (i.e., 12811 images) following (Zhai et al., 2019). The input is pre-processed using random crops resized to 224 x 224 followed by horizontal flip (p=0.5) and normalized to zero mean and one standard deviation using statistics of initial 10% partition.

**Architecture**: We use VGG16 (Simonyan & Zisserman, 2014) with batchnorm (Ioffe & Szegedy, 2015), 18-layer ResNet (He et al., 2016), and 28-layer 2-head Wide-ResNet (WRN-28-2) (Zagoruyko & Komodakis, 2016) in our experiments. For both target architectures we use<sup>3,4</sup>. For CIFAR10/100 models we set the number of neurons in penultimate fully-connected layer of VGG16 to 512 as in <sup>4</sup>.

**Regularization Hyper-parameters:** CIFAR10,  $lr=5e^{-4}$ ;wd=0 and CIFAR100,  $lr=5e^{-5}$ ;wd=0. Adam optimizer is used for both datasets. RA parameters are: CIFAR10: n=1, m=5, CIFAR100: n=1, m=2, ImageNet: n=2, m=9. We empirically select the SWA hyperparameters as: CIFAR 10/100: SWA LR: $5e^{-4}$  and frequency:50. Imagenet: SWA LR: $1e^{-5}$  and frequency:50. These parameters are selected after performing a grid search and kept consistent across experiments. We always train a model from scratch in each AL iteration except for Imagenet due to its heavy compute budget.

**Implementation of AL methods**: We developed a PyTorch-based toolkit to evaluate the AL methods in a unified implementation. AL methods can be cast into two categories based on whether or not AL sampling relies on the task model (classifier network). For example, coreset uses the latent space representations learnt by task model to select the sample set, whereas VAAL relies on a separate

VAE-discriminator network to select the samples, independent of the task model. In our implementation, we abstract these two approaches in a sampling function that may use the task model if required by the AL method. Each AL method was implemented using a separate sampling function, by referencing author-provided code if it was available. Using command line arguments, the toolkit allows the user to configure various aspects of training such as architecture used for task model, AL method, size of initial labeled set, size of acquisition batch, number of AL iterations, hyperparameters for task model training and AL sampling and number of repetitions.

## 5. Experiments and Results

All experiments were performed using 2 available nVidia DGX-1 servers, with each experiment utilizing 1–4 GPUs out of available 8 GPUs on each server. All codes were written in Python using PyTorch and other libraries in addition to third-party codebases. We plan to release our codebase on GitHub soon, for early-access please contact the authors.

#### 5.1. Variance in Evaluation Metrics

Training a neural network involves many stochastic components including parameter initialization, data augmentation, mini-batch selection, and batchnorm whose parameters change with mini-batch statistics. These elements can lead to a different optima thus resulting in varying performances across different runs of the same experiment. To evaluate the variance in classification accuracy caused by different initial labeled data, we draw five random initial labeled sets  $(L_0 \dots L_4)$  with replacement. Each of these five sets were used to train the base model, initialized with random weights, 5 times; a total of 25 models were trained for each AL method to characterize variance within-sample-sets and between-sample-sets.

From the results summarized in Fig. 1, we make the following observations: (i) A standard deviation of 1 to 2.5% in accuracy among different AL methods, indicating that out of chance, it is possible to achieve seemingly better results. (ii) In contrast to previous studies, our extensive experiments indicate that compared to RSB, no AL method achieves strictly better classification accuracy. At times, RSB appears to perform marginally better; for example, it achieves best mean accuracy of 80.36% (on CIFAR10 with 30% labeled data) and 35.72% (on CIFAR100 with 20%labeled data), whereas the second best performance is given by DBAL and VAAL i.e., 80.25\% and 35.59\% respectively. (iii) Our results averaged over 25 runs in Fig. 1 (f) and (l) indicate that no method performs clearly better than others. An ANOVA and pairwise multiple comparisons test with Tukey-Cramer FWER correction revealed that no AL method's performance was significantly different from RSB.

https://github.com/meliketoy/wide-resnet.pytorch
https://github.com/kuangliu/pytorch-cifar

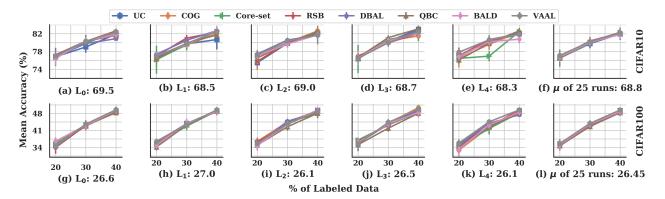


Figure 1. Comparisons of AL methods on CIFAR10 (top) and CIFAR100 (bottom) for different initial labeled sets  $L_0, L_1, \cdots, L_4$ . The mean accuracy for the base model (at 10% labeled data) is noted at the bottom of each subplot. The model is trained 5 times for different random initialization seeds. The mean of 25 runs in (f) & (l) suggest that no AL method performs significantly better than others. For exact numbers used to create the plots above, please refer to tables in the supplementary section.

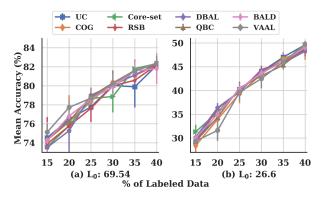


Figure 2. Results when only 5% of training data is annotated at each iteration of AL on (a) CIFAR10 and (b) CIFAR100. Results are average of 5 runs. For exact numbers used to create the plots above, please refer to tables in the supplementary section.

This provides a strong evidence and need to repeat an experiment over multiple runs to demonstrate true effectiveness of an AL method.

## 5.2. Differing Experimental Conditions

Next, we compare AL methods and RSB by modifying different experimental conditions for annotation batch size, size of validation set and class imbalance.

Annotation Batch Size (b): Following previous studies, we experiment with annotation batch size (b) equal to 5%, and 10% of the overall sample count (L+U). Results in Fig. 2 (corresponding table in supplementary section) show that VAAL and UC perform marginally better than the RSB, although this is inconsistent. For example, on CIFAR100 at 20% labeled data, and b=10%, VAAL performs marginally better than most of the AL methods (Fig. 1(1)). This is in contrast to results with b=5% (Fig. 2). We therefore conclude that no AL method offers consistent advantage

over others under different budget size settings.

Validation Set Size: During training, we select the best performing model on the validation set (V) to report the test set  $(T_s)$  results. To evaluate the sensitivity of AL results to the size of V, we perform experiments on CIFAR100 with three different V sizes: 2%, 5%, and 10% of the total samples (L+U). From results in Table 1, we do not observe any appreciable trend in accuracy with respect to the size of V. For example, the RSB achieves a mean accuracy of 49.8%, 49.1%, and 48.4%, respectively, for the best model selected using 2%, 5% and 10% of the training data as V. We conclude that AL results do not change significantly with the size of V, and a small V set can work for model selection in low-data regimes such as AL, freeing up more data for training the task model; a similar observation was made in a recent SSL study (Zhai et al., 2019).

Class Imbalance: Here, we evaluate the robustness of different AL methods on imbalanced data. For this, we construct  $L_0$  on CIFAR100 dataset, to simulate long tailed distribution of classes by following a power law, where the number of samples of 100 classes are given by samples [i] =  $a + b * \exp^{\alpha x}$  where  $i \in \{1...100\}; a = 100, x = 100$  $i + 0.5, \alpha = -0.046$  and b = 400. The resulting sample count per class is normalized to construct a probability distribution. Models were trained using previously described settings, with the exception of loss function which was set to weighted cross entropy. The results in Fig. 4 show that for the first two AL iterations, RSB achieves the highest mean accuracy (n = 5), and is surpassed by DBAL in the last iteration. More importantly, we notice that AL methods demonstrate different degree of change in the imbalanced class setting, without revealing a clear trend in the plot. In contrast to the previously reported observations that found AL methods robust to class imbalance in the dataset, we conclude that AL methods do not outperform RSB.

		2%			5%			10%	
Methods	20%	30%	40%	20%	30%	40%	20%	30%	40%
RSB	$34.6 \pm 1.2$	$43.3 \pm 1.6$	$\textbf{49.8} \pm \textbf{1.1}$	$35.4 \pm 1.4$	$42.5 \pm 1.9$	$49.1 \pm 1.7$	$34 \pm 0.3$	$43.1 \pm 1.5$	$48.4 \pm 1.1$
VAAL	$34.9 \pm 0.8$	$43.9 \pm 0.1$	$48.6 \pm 0.9$	$34.9 \pm 0.5$	$42.9\pm1.3$	$47.7\pm1.4$	$34.6 \pm 0.5$	$\textbf{43.6} \pm \textbf{0.8}$	$\textbf{49.5} \pm \textbf{0.9}$
UC	$36.8 \pm 0.7$	$43.7 \pm 0.5$	$48.8\pm1.2$	$33.7 \pm 2.2$	$43.5\pm1.2$	$49.1\pm0.4$	$34.9 \pm 1.1$	$42.8\pm1.7$	$48.9 \pm 0.7$
Coreset	$36.2 \pm 1.1$	$42.8\pm1.3$	$49.1 \pm 1.1$	$34.5 \pm 1.7$	$\textbf{44.4} \pm \textbf{0.7}$	$49.3 \pm 1.3$	$35.5 \pm 0.8$	$43.2 \pm 0.7$	$48.8 \pm 0.6$
COG	$35.4 \pm 1.4$	$\textbf{44.2} \pm \textbf{0.9}$	$49.2 \pm 1$	$34.1 \pm 2.1$	$43.7 \pm 0.7$	$48.8 \pm 1.7$	$35.9 \pm 2.2$	$42.7 \pm 1.4$	$49.4 \pm 1.4$
DBAL	$35.0 \pm 0.8$	$43.8\pm1.3$	$48.5\pm1.6$	$36.4 \pm 1.5$	$42.8 \pm 0.7$	$\textbf{50.0} \pm \textbf{0.8}$	$34.2 \pm 1.7$	$43.4 \pm 1.8$	$49.3 \pm 0.9$
BALD	$34.1 \pm 1.3$	$44 \pm 1$	$49.4 \pm 1$	$36.2 \pm 1.3$	$42.2\pm1.2$	$48.5 \pm 0.6$	$36.5 \pm 1.2$	$43.1 \pm 0.9$	$49.3 \pm 0.6$
QBC	$35.3 \pm 1.8$	$43.3 \pm 0.4$	$48.7\pm1$	$34.2 \pm 0.9$	$43.1\pm1.1$	$48.4 \pm 0.9$	$34.7 \pm 2.2$	$43.1\pm1.6$	$48.3 \pm 0.6$

Table 1. Test set performance for model selected with different validation set sizes on CIFAR100. Results are average of 5 runs.

Methods	CIFAR10	CIFAR100
RSB	$69.54 \pm 1.58$	$26.58 \pm 0.29$
+ SWA	$74.57 \pm 0.87$	$32.51 \pm 0.92$
+ RA	$75.43 \pm 0.89$	$29.77 \pm 0.83$
+ Shake-Shake(SS)	$71.78 \pm 0.99$	$34.8 \pm 0.28$
+ SWA $+$ RA	$79.86 \pm 0.6$	$36.65 \pm 0.35$
+ SS + SWA + RA	$\textbf{82.88} \pm \textbf{0.26}$	$44.37 \pm 0.78$

Table 2. Individual Contributions of different regularization techniques. Results averaged over 5 runs for 10% of training data. Above all experiments use the VGG16 architecture except for Shake-Shake as it is restricted to the family of resnext.

#### 5.3. Regularization

With the motivation stated in section 3, we evaluate the effectiveness of advanced regularization techniques (RA and SWA) in the context of AL using CIFAR10 and CIFAR100 datasets. All experimental settings were used as previously reported, with the exception of number of epochs which was increased to 150 (from 100). We empirically observed that unlike  $\ell_2$ —regularization, which requires careful tuning, RA and SWA work fairly well with changes in their hyperparameters. We therefore do not use  $\ell_2$ —regularization in these experiments where RA and SWA was applied.

Fig. 3 compares different AL methods with RSB on CI-FAR10/100 datasets. We observe that models trained with RA and SWA consistently achieve significant performance gains across all AL iterations and exhibit appreciablysmaller variance across multiple runs of the experiments. Our regularized random-sampling baselines on 40% labeled data achieves mean accuracy of 89.73% and 57.16% respectively on CIFAR10 and CIFAR100. We note that using RSB, for CIFAR10, a model regularized using RA and SWA with 20% of training data achieves over 4% higher accuracy compared to a model trained without RA and SWA using much larger 40% of the training data. Similarly for CIFAR100, the RSB 20%-model with regularization performs comparably to the 40%-model without regularization. Therefore, we consider regularization to be a valuable addition to the low-data training regime of AL, especially given

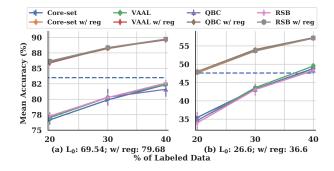


Figure 3. Effect of regularization (RA + SWA) on the test accuracy of CIFAR10(a) and CIFAR100(b) dataset. Results are average of 5 runs where regularized results are shown above the blue line.

that it significantly reduces the variance in evaluation metric and helps avoid misleading conclusions.

An ablative study to show individual contribution of each regularization technique towards overall performance gain is given in Table 2. The results indicate that both RA and SWA show a significant combined gain of  $\approx 10\%$ . We also experimented with Shake-Shake (SS) (Gastaldi, 2017) in parallel to RA and SWA, and observed that it significantly increases the runtime, and is not robust to model architectures. We therefore chose RA & SWA over SS in our experiments.

### 5.4. Transferability and Optimizer Settings

In principle, the sample sets drawn by an AL method should be agnostic to the task model's architecture, and a change in the architecture should maintain consistent performance trends for the AL method. We conduct an experiment by storing the indices of sample set drawn in an AL iteration on the source network, and use them to train the target network. We consider VGG16 as the source, and ResNet18 (RN18) (He et al., 2016) & WRN-28-2 (Zagoruyko & Komodakis, 2016) as the target architectures. From Table 3, we observe that the trend in AL gains is architecture dependent. On CIFAR10 with RN18 using Adam, VAAL achieves higher accuracy than RSB. However, this relative gain vanishes with RA and SWA. Further, there was no discernible trend in results using WRN-28-2 or VGG16 architectures.

	Sou	ırce Mo	del		Target Model													
					R18			R18 R18					R18					
		VGG16	)	W	WRN-28-2 +Adam			+SGD +		+4	+Adam+Reg		+SGD+Reg					
	20%	30%	40%	20%	30%	40%	20%	30%	40%	20%	30%	40%	20%	30%	40%	20%	30%	40%
RSB	77.3	80.3	82.6	79.1	82.4	84.7	74.1	77.3	80.8	80.1	84.1	86.2	86.7	89	90.4	84.8	87.8	89.3
Coreset	76.7	79.9	82.4	79.1	82.9	83.7	74.4	78.8	81.1	80.1	84	86.5	86.4	88.9	90.3	85.1	87.2	89.2
VAAL	77.0	80.3	82.4	78.9	82.7	84.1	75.7	<b>79.6</b>	81.5	79.6	83.8	86.4	86.6	88.9	90.5	84.9	87.7	89.3
QBC	77.2	80.3	81.6	78.1	82.9	84.9	74.3	77.8	80.6	79.9	83.6	86.1	86.6	88.9	90.1	85.1	87.6	89.3

Table 3. Transferability experiment on CIFAR10 dataset where source model is VGG16 and target model is Resnet18 (R18) and Wide Resnet-28-2 (WRN-28-2). The reported numbers are mean of test accuracies over five seeds on CIFAR10/100 dataset. Results with regularization (Reg=SWA+RA) are shown in last two columns.

To evaluate whether the choice of optimizer played a role in VAAL's performance using RN18 with Adam, we repeated the training with SGD. We note the followings (Table 3): (i) RSB (and other methods) achieved a higher mean accuracy when trained using SGD compared to Adam (74.1% vs 80.1%) on RN18 using 20% CIFAR10 labeled data. Further, RN18 with SGD performs comparably against WRN-28-2 with Adam *i.e.*, 80.1% vs 79.1%. (ii) Using Adam, both VAAL and coreset perform favorably against RSB. However, with SGD, the results are comparable.

## 5.5. Active Learning on ImageNet

Compared to CIFAR10/100, ImageNet is more challenging with larger sample count, 1000 classes and higher resolution images. We compare coreset, VAAL and RSB on ImageNet. We were unable not evaluate QBC due to prohibitive compute cost of training an ensemble of 5 CNN models. The details for training hyper-parameters are in supplementary section. Results with and without regularization (RA, SWA) are shown in Table 5. Using ResNext-50 architecture (Xie et al., 2017) and following the settings of (Zhai et al., 2019)), we achieve improved baseline performances compared to the previously reported results (Beluch et al., 2018; Sinha et al., 2019). From table 5, we observe that both AL methods performed marginally better than RSB though ImageNet experiments are not repeated for multiple runs due to prohibitive compute requirements.

#### 5.6. Additional Experiments

Noisy Oracle: In this experiment, we sought to evaluate the stability of regularized network to labels from a noisy oracle. We experimented with two levels of oracle noise by randomly permuting labels of 10% and 20% of samples in the set drawn by random sampling baseline at each iteration. From results in Table 4, we found that the drop in accuracy for the model regularized by RA and SWA was nearly half (3%) compared with the model trained without these regularizations (6%) on both 30% and 40% data splits. Our findings suggest that the noisy pseudolabels generated for the unlabelled set U by model  $\phi$ , when applied in conjunction with appropriate regularization, should help improve model's performance. Additional results using AL meth-

ods in this setting are shared in the supplementary section. **Active Learning Sample Set Overlap**: For interested readers, we discuss the extent of overlap among the sample sets drawn by AL methods in the supplementary section.

#### 6. Discussion

Under-Reported Baselines: We note that several recent AL studies show baseline results that are lower than the ones reproduced in this study. Table 6 summarizes our RSB results with comparisons to some of the recently published AL methods, under similar training settings. Based on this observation, we emphasize that comparison of AL methods must be done under a consistent set of experimental settings. Our observations confirm and provide a stronger evidence for a similar conclusion drawn in (Mittal et al., 2019), and to a less related extent, (Oliver et al., 2018). Different from (Mittal et al., 2019) though, we demonstrate that: (i) relative gains using AL method are found under a narrow combination of experimental conditions, (ii) such gains are not statistically meaningful over random baseline, (iii) more distinctly, we show that the performance gains vanish when a well-regularized training strategy is used.

The Role of Regularization: Regularization helps reduce generalization error and is particularly useful in training overparameterized neural networks with low data. We show that both RA and SWA can achieve appreciable gain in performance at the expense of a small computational overhead. We observed that along with learning rate (in case of SGD), regularization was one of the key factors in reducing the error while being fairly robust to its hyperparameters (in case of RA and SWA). We also found that any trend of consistent gain observed with an AL method over RSB on CIFAR10/100 disappears when the model is wellregularized. Models regularized with RA and SWA also exhibited smaller variance in evaluation metric compared to the models trained without them. With these observations, we recommend that AL methods be also tested using well-regularized model to ensure their robustness. Lastly, we note that there are multiple ways to regularize the datamodel-metric pipeline, we focus on data and model side regularization using techniques such as RA and SWA, though it is likely that other combination of newer regularization

Methods ↓	10%	20%	30%	40%					
Noise: 10%									
RSB	69.09	72.78	76.97	76.63					
RSB + Reg.	79.28	85.02	87.05	88.01					
	Nois	e: 20%							
RSB	69.09	70.37	71.01	70.04					
RSB + Reg.	79.28	83.02	84.24	85.44					

Table 4. RSB accuracy with and without SWA and RA on CIFAR10 with noisy oracle. RSB+Reg. refers to RSB regularized using RA and SWA.

$Methods \downarrow$	10%	15%	20%	25%				
without RA + SWA								
RSB	58.05	62.95	64.61	66.15				
VAAL	58.05	63.33	64.68	66.18				
Coreset	58.05	63.04	64.43	65.58				
	with F	A + SW	A					
RSB	59.43	63.88	66.83	69.10				
VAAL	59.43	65.17	67.39	69.47				
Coreset	59.43	64.17	67.07	69.54				

Table 5. Effect of RA and SWA on ImageNet where annotation budget is 5% of training data. Results reported for 1 run.

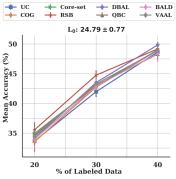


Figure 4. Results are average of 5 runs on imbalanced CIFAR100.

techniques will lead to similar results. We do believe that with their simplicity and applicability to a wide variety of model (as compared to methods such as shake-shake), RA and SWA can be effectively used in AL studies without significant hyperparameter tuning.

Methods	10%	20%	30%	40%					
	CIFAR10								
VAAL	61.35	68.17	72.26	75.99					
Coreset	60	68	71	74					
RSB(ours)	69.54	77.29	80.28	82.61					
RSB-R(ours)	79.86	86.18	88.36	89.73					
CIFAR100									
VAAL	28.8	35.35	41.7	45.9					
Coreset	29	37	42	48					
RSB(ours)	26.58	33.99	43.08	48.38					
RSB-R(ours)	36.65	43.89	50.07	54.58					

*Table 6.* Reported Random Baseline Accuracies vs our RSB results. We denote our RSB results with regularization by RSB-R.

Using Unlabeled Set in Training: Some recent methods such as VAAL use U set to train another network as part of their sampling routine. We argue that for such models, a better baseline comparison would be from the semi-supervised learning (SSL) literature. We note that some of the current SSL methods such as UDA (Xie et al., 2019) have reported very strong results (94.71% on CIFAR10 with 8% labeled training data). These results suggest that large number of noisy labels are relatively more helpful in reducing the generalization error as compared to the smaller percentage of high quality labels. Further commentary on this topic can be found in (Mittal et al., 2019).

AL Methods Compared To Strong RSB: Compared to the well-regularized RSB, state-of-the-art AL methods evaluated in this paper do not achieve any noticeable gain. We believe that reported AL results in the literature were obtained with insufficiently-regularized models, and the gains

reported for AL methods are often not because of the superior quality of selected samples. As shown in Table 3, the fact that a change in model architecture can change the conclusions being drawn suggests that transferability experiments should be essential to any AL study. Similarly we observed that a simple change in optimizer or use of regularization can influence the conclusions. The highly-sensitive nature of AL results using neural networks therefore necessitates a comprehensive suite of experimental tests.

# 7. Conclusion and Proposed Guidelines

Our extensive experiments suggest a strong need for a common evaluation platform that facilitates robust and reproducible development of AL methods. To this end, we recommend the following to ensure results are robust: (i) experiments must be repeated under varying training settings such as optimizer, and model architecture, budget size, among others, (ii) regularization techniques such as RA and SWA should be incorporated into the training to ensure AL methods are able to demonstrate gains over a regularized random baseline, (iii) transferability experiments must be performed to ensure the AL-drawn sample sets are indeed informative as claimed. To increase the reproducibility of AL results, we further recommend: (iv) experiments should be performed using a common evaluation platform under consistent settings to minimize the sources of variation in the evaluation metric, (v) snapshot of experimental settings should be shared, e.g. using a configuration file (.cfg, .json etc), (vi) index sets for a public dataset used for partitioning the data into training, validation, test, and AL-drawn sets should be shared, along with the training scripts. In order to facilitate the use of these guidelines in AL experiments, we also provide an open-source AL toolkit. We believe our findings and toolkit will help advance robust and reproducible AL research.

#### References

- Beluch, W. H., Genewein, T., Nürnberger, A., and Köhler, J. M. The power of ensembles for active learning in image classification. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 9368–9377, 2018.
- Cubuk, E. D., Zoph, B., Mane, D., Vasudevan, V., and Le, Q. V. Autoaugment: Learning augmentation policies from data. *arXiv preprint arXiv:1805.09501*, 2018.
- Cubuk, E. D., Zoph, B., Shlens, J., and Le, Q. V. Randaugment: Practical data augmentation with no separate search. *arXiv* preprint arXiv:1909.13719, 2019.
- Ducoffe, M. and Precioso, F. Adversarial active learning for deep networks: a margin based approach. *CoRR*, abs/1802.09841, 2018. URL http://arxiv.org/abs/1802.09841.
- Gal, Y., Islam, R., and Ghahramani, Z. Deep bayesian active learning with image data. In *Proceedings of the 34th International Conference on Machine Learning-Volume* 70, pp. 1183–1192. JMLR. org, 2017.
- Gastaldi, X. Shake-shake regularization. *arXiv preprint arXiv:1705.07485*, 2017.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B.,
  Warde-Farley, D., Ozair, S., Courville, A., and Bengio,
  Y. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pp. 2672–2680, 2014.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- Izmailov, P., Podoprikhin, D., Garipov, T., Vetrov, D., and Wilson, A. G. Averaging weights leads to wider optima and better generalization. *arXiv preprint* arXiv:1803.05407, 2018.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015.
- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Kirsch, A., van Amersfoort, J., and Gal, Y. Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning. *CoRR*, abs/1906.08158, 2019. URL http://arxiv.org/abs/1906.08158.

- Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- Lewis, D. D. and Gale, W. A. A sequential algorithm for training text classifiers. In SIGIR94, pp. 3–12. Springer, 1994.
- Lowell, D., Lipton, Z. C., and Wallace, B. C. How transferable are the datasets collected by active learners? *CoRR*, abs/1807.04801, 2018. URL http://arxiv.org/abs/1807.04801.
- Mittal, S., Tatarchenko, M., zgn iek, and Brox, T. Parting with illusions about deep active learning, 2019.
- Oliver, A., Odena, A., Raffel, C. A., Cubuk, E. D., and Goodfellow, I. Realistic evaluation of deep semi-supervised learning algorithms. In *Advances in Neural Information Processing Systems*, pp. 3235–3246, 2018.
- Prabhu, A., Dognin, C., and Singh, M. Sampling bias in deep active classification: An empirical study. *arXiv* preprint arXiv:1909.09389, 2019.
- Sener, O. and Savarese, S. Active learning for convolutional neural networks: A core-set approach. In *International Conference on Learning Representations*, 2018. URL https://openreview.net/forum?id=H1aIuk-RW.
- Simonyan, K. and Zisserman, A. Very deep convolutional networks for large-scale image recognition. *arXiv* preprint arXiv:1409.1556, 2014.
- Sinha, S., Ebrahimi, S., and Darrell, T. Variational adversarial active learning. *arXiv preprint arXiv:1904.00370*, 2019.
- Tran, T., Do, T., Reid, I. D., and Carneiro, G. Bayesian generative active deep learning. *CoRR*, abs/1904.11643, 2019. URL http://arxiv.org/abs/1904.11643.
- Xie, Q., Dai, Z., Hovy, E., Luong, M.-T., and Le, Q. V. Unsupervised data augmentation for consistency training. *arXiv* preprint arXiv:1904.12848, 2019.
- Xie, S., Girshick, R., Dollár, P., Tu, Z., and He, K. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1492–1500, 2017.
- Yoo, D. and Kweon, I. S. Learning loss for active learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 93–102, 2019.
- Zagoruyko, S. and Komodakis, N. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.

Zhai, X., Oliver, A., Kolesnikov, A., and Beyer, L.  $S^4$ l: Self-supervised semi-supervised learning. *CoRR*, abs/1905.03670, 2019. URL http://arxiv.org/abs/1905.03670.

# **Supplementary Section**

## 1. Hyper-parameters

In this section we mention the training details which were used to report the experiments in main paper.

## 1.1. Transferability Experiment

We mainly used three different architectures for classifier model *i.e.*VGG16, ResNet18 (R18) and Wide ResNet-28-2 (WRN)<sup>1</sup>. The VGG network was used as a source model whereas other two networks are used for target models.

- (i) VGG16  $\longrightarrow$  R18
  - (a) SGD Optimizer: 200 epochs; lr = 0.1 (decays by a factor of 10 at epoch steps: 100 and 150)
  - (b) Adam Optimizer: When RA and SWA was used, we trained for 115 epochs, otherwise 100 epochs;  $lr = 5e^{-4}$ , wd = 0.
  - (c) Regularization details are reported in the following table:

		SWA		RA		
Optimizer	LR	Frequency	Epochs	#Transforms	Index Magnitude	
SGD	$1e^{-3}$	50	35	1	5	
Adam	$5e^{-4}$	50	35	1	5	

Table 1. Regularization Hyper-parameters

- (ii) VGG16 → WRN-28-2: Some details
  - (a) Adam Optimizer: 100 epochs;  $lr = 5e^{-4}$ ,  $wd = 5e^{-4}$
  - (b) Results for CIFAR100 are reported in Table 2 which are achieved when we replace all the relu activations with leaky relu (negative slope set to 0.2) following (Oliver *et al.*, 2018). We found CIFAR100 results to be significantly better with leaky relu activation, however, the same change does not affect the performance of CIFAR10.

#### 2. Overlap in Active-set

Results are summarized in Figure 1.

## 3. Annotation Batch Size

Here we present the results for CIFAR10 and CIFAR100 in Table 15 for the experiment where annotation batch size is 5% relative to training data.

#### 4. Noisy Oracle Experiments

In conjunction to RSB baselines (presented in main paper), we report performance of AL methods under noisy labels in active sets. The results are reported in table 4 where we make the following observations: (i) it is quite evident that both SWA and RA improves performance even when label corruptions scenarios. (ii) No AL method consistently outperforms the simple RSB baseline. (iii) SWA and RA help reduce the performance difference between RSB and best AL method at a particular data split.

<sup>&</sup>lt;sup>1</sup>All Model definitions have been provided in AL toolkit.

	Sou	arce Mo	del	Target Model					
	VGG16			W	/RN-28	-2	R18+Adam		
Methods $\downarrow$	20%	30%	40%	20%	30%	40%	20%	30%	40%
Random	34.0	43.1	48.4	47.3	54.7	58.8	46.3	53.0	57.1
Coreset	35.5	43.2	48.8	48.6	54.3	58.4	47.3	52.9	<b>57.</b> 5
VAAL	34.6	43.6	49.5	48.0	54.3	58.4	46.5	52.9	57.1
QBC	34.7	43.1	48.3	48.4	54.0	58.5	46.2	53.0	57.3

Table 2. Transferability experiment on CIFAR100 dataset where source model is VGG16. The reported numbers are mean of test accuracies over 5 seeds. For this experiment we replace all relu activations with leaky relu (negative slope set to 0.2).

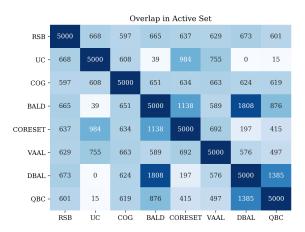


Figure 1. Overlap between active sets computed on CIFAR100 for AL iteration of 10% to 20%.

# 5. ImageNet Results

Our first run with ImageNet shown in the main paper suggested that both AL methods (Coreset and VAAL) were better than RSB across the three AL iterations when RA and SWA was used. In order to verify if this trend holds, we repeated ImageNet AL experiment to verify if the findings are repeatable. The table 3 below shows the results of both the runs. At 15%, RSB achieves higher accuracy compared to other two methods, and at 20% it performs better than VAAL but marginally lower than Coreset. Based on these results, we conclude that AL methods do not offer a consistent advantage over RSB.

### 6. Additional Results

In the last we present the exact accuracies which were used to plot the Figure 1 in main paper.

Methods ↓	10%	15%	20%	25%						
Run 1										
RSB	59.43	63.88	66.83	69.10						
VAAL	59.43	65.17	67.39	69.47						
Coreset	59.43	64.17	67.07	69.54						
	Run 2									
RSB	59.28	64.98	67.31	-						
VAAL	59.28	64.28	66.89	-						
Coreset	59.28	64.22	67.48	-						

Table 3. Effect of RA and SWA on Imagenet where annotation budget is 5% of training data.

	wi	thout re	gulariza	tion	W	ith regu	larizatio	n	
<b>Methods</b> ↓	10%	20%	30%	40%	10%	20%	30%	40%	
Noise: 10%									
UC	68.5	69.11	70.89	72.39	79.95	81.54	81.91	82.87	
COG	68.5	68.59	70.02	72.46	79.95	81.78	82.29	83.04	
Coreset	68.5	72.89	76	78.05	79.95	84.97	87.39	88.16	
RSB	68.5	73.6	77.02	78.49	79.95	85.37	87.46	88.59	
DBAL	68.5	73.06	76.83	79.47	79.95	85.32	87.92	89.34	
QBC	68.5	73.44	76.4	78.7	79.95	85.28	87.66	89.04	
BALD	68.5	71.5	74.7	76.92	79.95	84.98	87.29	88.96	
VAAL	68.5	72.93	75.52	76.76	79.95	84.01	85.27	86.58	
			Noi	se: 20%					
UC	68.5	67.15	66.94	69.28	79.95	81.04	80.45	79.54	
COG	68.5	66.74	66.7	68.29	79.95	80.9	80.43	78.59	
Coreset	68.5	69.77	71.64	73.44	79.95	82.92	84.84	85.74	
RSB	68.5	70.44	71.43	72.87	79.95	83.75	85.22	86.94	
DBAL	68.5	71.24	70.99	73.21	79.95	84.01	85.27	86.66	
QBC	68.5	69.63	69.44	72.11	79.95	83.68	85.36	86.35	
BALD	68.5	70.44	71.86	74.06	79.95	83.35	85.31	85.89	
VAAL	68.5	69.09	69.19	71.65	79.95	81.94	82.91	83.26	

Table 4. Mean accuracy on noisy oracle experiments on CIFAR10 with (n=3) trials. We note that the noise is added in active sets drawn by AL methods. The regularization experiments involve SWA and RA techniques.

Methods	20%	30%	40%
RSB	$77.29 \pm 0.38$	$80.28 \pm 0.70$	$\textbf{82.61} \pm \textbf{0.50}$
UC	$76.82 \pm 1.48$	$79.02 \pm 1.30$	$81.89 \pm 0.67$
COG	$76.96 \pm 1.26$	$80.37 \pm 0.72$	$81.92 \pm 0.70$
Coreset	$76.66 \pm 0.76$	$79.90 \pm 0.65$	$82.37 \pm 0.63$
VAAL	$77.02 \pm 0.87$	$80.29 \pm 0.99$	$82.40 \pm 0.58$
DBAL	$76.80 \pm 2.03$	$79.76 \pm 1.32$	$80.88 \pm 0.36$
QBC	$77.24 \pm 0.50$	$80.33 \pm 1.36$	$81.61\pm1.20$
BALD	$76.46 \pm 1.49$	$\textbf{80.41} \pm \textbf{0.66}$	$81.80 \pm 0.99$

Methods	20%	30%	40%
RSB	$77.02 \pm 0.68$	$\textbf{81.01} \pm \textbf{0.71}$	$82.31 \pm 0.42$
UC	$76.75 \pm 1.21$	$79.69\pm1.32$	$80.69 \pm 2.25$
COG	$76.36 \pm 0.70$	$79.45\pm0.80$	$82.03 \pm 0.63$
Coreset	$76.17 \pm 3.14$	$79.57 \pm 1.34$	$82.51 \pm 1.1$
VAAL	$76.90 \pm 1.12$	$79.8 \pm 0.76$	$82.51 \pm 0.42$
DBAL	$\textbf{77.38} \pm \textbf{0.59}$	$80.53 \pm 0.55$	$\textbf{82.66} \pm \textbf{0.80}$
QBC	$76.33 \pm 1.03$	$79.92 \pm 1.28$	$81.70 \pm 1.20$
BALD	$77.07 \pm 0.54$	$79.64 \pm 0.76$	$82.70 \pm 0.69$

Table 5. CIFAR10 Test Accuracy on  $L_0^0$ . The base model accuracy is  $69.54 \pm 1.58$ . Table 6. CIFAR10 Test Accuracy on  $L_0^1$ . The base model accuracy is  $68.55 \pm 2.91$ .

Methods	20%	30%	40%	
RSB	$75.80 \pm 0.95$	$79.89 \pm 0.81$	$81.86 \pm 0.40$	
UC	$75.50 \pm 1.38$	$79.85 \pm 0.97$	$82.16\pm1.13$	
COG	$75.78 \pm 1.94$	$80.00 \pm 0.47$	$\textbf{82.62} \pm \textbf{0.23}$	
Coreset	$77.04 \pm 0.65$	$79.92 \pm 0.65$	$82.31 \pm 0.96$	
VAAL	$77.27 \pm 1.07$	$80.40 \pm 0.30$	$81.65 \pm 2.05$	
DBAL	$77.45 \pm 0.85$	$\textbf{80.57} \pm \textbf{0.47}$	$81.67 \pm 1.06$	
QBC	$76.54 \pm 1.03$	$79.84 \pm 0.72$	$82.49 \pm 0.53$	
BALD	$77.03 \pm 0.78$	$79.69 \pm 0.77$	$81.75 \pm 0.49$	

Methods	20%	30%	40%	
RSB	$77.01 \pm 0.75$	$79.95 \pm 0.85$	$82.05 \pm 0.81$	
UC	$76.39 \pm 1.23$	$80.26\pm0.73$	$82.98 \pm 0.29$	
COG	$76.85 \pm 1.09$	$80.41 \pm 0.50$	$81.49 \pm 0.97$	
Coreset	$76.31 \pm 3.15$	$80.67\pm0.58$	$82.01 \pm 1.72$	
VAAL	$76.47 \pm 1.15$	$79.98 \pm 0.62$	$82.38 \pm 0.56$	
DBAL	$76.63 \pm 1.10$	$80.15\pm0.96$	$82.67 \pm 0.39$	
QBC	$76.66 \pm 0.91$	$\textbf{81.11} \pm \textbf{0.47}$	$\textbf{83.08} \pm \textbf{0.71}$	
BALD	$76.55 \pm 1.55$	$80.37\pm0.63$	$82.06\pm0.98$	

Table 7. CIFAR10 Test Accuracy on  $L_0^2$ . The base model accuracy Table 8. CIFAR10 Test Accuracy on  $L_0^3$ . The base model accuracy is  $69.03 \pm 2.15$ .

is  $68.70 \pm 1.96$ .

Methods	20%	30%	40%	
RSB	$77.05 \pm 1.84$	$\textbf{80.67} \pm \textbf{0.30}$	$81.96 \pm 0.39$	
UC	$77.11 \pm 1.12$	$80.04\pm0.87$	$81.98 \pm 0.41$	
COG	$76.40 \pm 1.99$	$79.61 \pm 1.27$	$82.34 \pm 0.45$	
Coreset	$76.46 \pm 1.15$	$79.96 \pm 1.16$	$82.47 \pm 0.64$	
VAAL	$\textbf{77.85} \pm \textbf{0.76}$	$80.62 \pm 1.10$	$81.85 \pm 0.43$	
DBAL	$76.78 \pm 0.97$	$80.27 \pm 1.12$	$82.18 \pm 0.86$	
QBC	$76.10 \pm 0.86$	$79.89 \pm 0.85$	$\textbf{82.63} \pm \textbf{0.66}$	
BALD	$76.89 \pm 1.35$	$80.21\pm0.66$	$80.75\pm0.93$	

Table 9. CIFAR10 Test Accuracy on  $L_0^4$ . The base model accuracy is  $68.31 \pm 1.60$ .

Methods	20%	30%	40%	
RSB	$33.99 \pm 2.59$	$43.08 \pm 1.46$	$48.38 \pm 1.10$	
UC	$34.87 \pm 1.14$	$42.78 \pm 1.74$	$48.88 \pm 0.72$	
COG	$35.91 \pm 2.21$	$42.68 \pm 1.41$	$49.41 \pm 1.42$	
Coreset	$35.46 \pm 0.80$	$43.23 \pm 0.67$	$48.84 \pm 0.62$	
VAAL	$34.62 \pm 0.47$	$\textbf{43.61} \pm \textbf{0.79}$	$\textbf{49.50} \pm \textbf{0.87}$	
DBAL	$34.25 \pm 1.66$	$43.44\pm1.76$	$49.26\pm0.87$	
QBC	$34.68 \pm 2.16$	$43.07\pm1.57$	$48.30 \pm 0.63$	
BALD	$36.55 \pm 1.25$	$43.14 \pm 0.94$	$49.28\pm0.56$	

20%	30%	40%	
$\textbf{36.48} \pm \textbf{0.56}$	$43.53 \pm 1.27$	$48.93 \pm 0.80$	
$35.62\pm1.63$	$\textbf{43.93} \pm \textbf{1.79}$	$49.07\pm0.57$	
$34.97 \pm 1.67$	$43.39\pm1.22$	$49.24 \pm 0.49$	
$35.34 \pm 1.91$	$42.78 \pm 1.71$	$49.05 \pm 1.01$	
$36.33 \pm 0.91$	$43.40\pm0.45$	$\textbf{49.43} \pm \textbf{0.80}$	
$35.94 \pm 1.69$	$43.91 \pm 1.21$	$48.66 \pm 1.43$	
$34.20 \pm 1.47$	$43.25\pm0.29$	$48.93 \pm 1.66$	
$34.93 \pm 1.48$	$43.93 \pm 1.92$	$48.94 \pm 0.97$	
	$36.48 \pm 0.56$ $35.62 \pm 1.63$ $34.97 \pm 1.67$ $35.34 \pm 1.91$ $36.33 \pm 0.91$ $35.94 \pm 1.69$ $34.20 \pm 1.47$	36.48 $\pm$ 0.5643.53 $\pm$ 1.2735.62 $\pm$ 1.6343.93 $\pm$ 1.7934.97 $\pm$ 1.6743.39 $\pm$ 1.2235.34 $\pm$ 1.9142.78 $\pm$ 1.7136.33 $\pm$ 0.9143.40 $\pm$ 0.4535.94 $\pm$ 1.6943.91 $\pm$ 1.2134.20 $\pm$ 1.4743.25 $\pm$ 0.29	

Table 10. CIFAR100 Test Accuracy on  $L_0^0$ . The base model accuracy is  $26.58 \pm 0.29$ . The base model accuracy is  $27.03 \pm 0.26$ .

Methods	20%	30%	40%	
RSB	$36.3 \pm 1.46$	$43.50 \pm 1.33$	$\textbf{49.41} \pm \textbf{0.94}$	
UC	$36.03 \pm 0.84$	$43.82 \pm 0.94$	$49.31 \pm 0.35$	
COG	$36.50 \pm 1.05$	$44.46 \pm 0.69$	$49.10 \pm 0.97$	
Coreset	$35.66 \pm 1.07$	$44.24\pm0.77$	$48.03 \pm 0.94$	
VAAL	$35.70 \pm 2.32$	$43.59 \pm 1.43$	$48.87 \pm 0.89$	
DBAL	$35.02 \pm 1.45$	$\textbf{44.80} \pm \textbf{0.33}$	$48.49 \pm 1.69$	
QBC	$34.97 \pm 1.53$	$42.51 \pm 1.62$	$48.07 \pm 2.04$	
BALD	$34.81 \pm 2.05$	$43.52 \pm 1.02$	$49.31 \pm 0.92$	

Methods	20%	30%	40%	
RSB	$36.45 \pm 0.34$	$43.31 \pm 1.01$	$49.16 \pm 0.87$	
UC	$35.85 \pm 0.98$	$43.85 \pm 0.71$	$49.32\pm1.06$	
COG	$35.77 \pm 1.13$	$\textbf{43.88} \pm \textbf{1.20}$	$\textbf{50.41} \pm \textbf{0.54}$	
Coreset	$\textbf{37.05} \pm \textbf{1.02}$	$43.27\pm0.64$	$48.4 \pm 1.98$	
VAAL	$35.47 \pm 2.00$	$44.30\pm1.44$	$49.77 \pm 1.64$	
DBAL	$35.73 \pm 1.32$	$43.53 \pm 1.51$	$49.01 \pm 1.08$	
QBC	$35.15 \pm 2.67$	$41.92 \pm 1.38$	$48.08 \pm 1.49$	
BALD	$36.07 \pm 1.77$	$43.77 \pm 0.97$	$48.4 \pm 1.71$	

Table 12. CIFAR100 Test Accuracy on  $L_0^2$ . The base model accuracy Table 13. CIFAR100 Test Accuracy on  $L_0^3$ . The base model accuracy racy is  $26.11 \pm 0.36$ .

is  $26.47 \pm 0.46$ .

Methods	20%	30%	40%	
RSB	$35.37 \pm 1.66$	$42.73 \pm 1.30$	$49.28 \pm 1.07$	
UC	$35.08 \pm 1.05$	$42.77 \pm 1.59$	$49.13 \pm 0.73$	
COG	$32.88 \pm 2.61$	$42.38 \pm 2.09$	$48.46 \pm 1.19$	
Coreset	$34.34 \pm 1.60$	$41.74 \pm 2.41$	$48.29\pm1.45$	
VAAL	$35.81 \pm 1.49$	$\textbf{44.55} \pm \textbf{1.05}$	$\textbf{49.44} \pm \textbf{0.78}$	
DBAL	$35.00 \pm 1.29$	$43.95\pm0.81$	$47.44 \pm 0.91$	
QBC	$34.37 \pm 2.64$	$42.36 \pm 0.92$	$48.23 \pm 0.99$	
BALD	$33.33 \pm 3.34$	$43.24\pm1.77$	$48.53 \pm 0.91$	

Table 14. CIFAR100 Test Accuracy on  $L_0^4$ . The base model accuracy is  $26.08 \pm 0.31$ .

Methods $\downarrow$	10%	15%	20%	25%	30%	35%	40%
				CIFAR10			
RSB	$69.54 \pm 1.58$	$73.90 \pm 2.81$	$75.77 \pm 1.52$	$77.64 \pm 1.46$	$79.99 \pm 1.08$	$80.58 \pm 0.9$	$82.17 \pm 1.03$
UC	$69.54 \pm 1.58$	$74.53 \pm 1.17$	$76.5 \pm 1.06$	$77.78 \pm 0.47$	$80.04 \pm 0.55$	$79.89 \pm 2.18$	$82.14 \pm 0.46$
COG	$69.54 \pm 1.58$	$74.29 \pm 0.96$	$76.2 \pm 2.82$	$78.35\pm1.02$	$\textbf{80.36} \pm \textbf{0.54}$	$81.07 \pm 0.5$	$82.17 \pm 0.67$
Coreset	$69.54 \pm 1.58$	$73.58 \pm 1.55$	$75.79 \pm 0.98$	$78.67\pm1.05$	$78.87 \pm 1.67$	$\textbf{81.81} \pm \textbf{0.95}$	$82.33 \pm 0.8$
VAAL	$69.54 \pm 1.58$	$\textbf{75.12} \pm \textbf{0.89}$	$\textbf{77.73} \pm \textbf{0.71}$	$78.7 \pm 1.11$	$80.22\pm1.07$	$81.57 \pm 0.31$	$\textbf{82.4} \pm \textbf{0.73}$
DBAL	$69.54 \pm 1.58$	$73.51 \pm 1.3$	$75.27 \pm 2.42$	$\textbf{79.01} \pm \textbf{0.29}$	$80.06 \pm 0.5$	$81.14 \pm 1.09$	$82.3 \pm 0.67$
QBC	$69.54 \pm 1.58$	$74.35\pm1.17$	$76.28\pm1.39$	$78.87 \pm 0.94$	$80.3 \pm 1.32$	$81.61 \pm 1.11$	$82.28 \pm 0.78$
BALD	$69.54 \pm 1.58$	$74.22\pm1.98$	$76.81\pm0.92$	$78.62 \pm 0.6$	$79.9 \pm 0.94$	$81.57\pm1.08$	$81.81 \pm 1.63$
				CIFAR100			
RSB	$26.58 \pm 0.29$	$29.43 \pm 1.36$	$36.22 \pm 0.89$	$39.70 \pm 1.20$	$\textbf{44.40} \pm \textbf{0.46}$	$45.88 \pm 0.89$	$48.71 \pm 1.43$
UC	$26.58 \pm 0.29$	$29.2\pm1.06$	$35.30 \pm 1.84$	$40.21\pm0.74$	$43.93 \pm 0.72$	$\textbf{47.03} \pm \textbf{0.94}$	$49.62 \pm 0.54$
COG	$26.58 \pm 0.29$	$28.42 \pm 1.33$	$33.92 \pm 1.41$	$39.33 \pm 2.01$	$43.45 \pm 1.27$	$46.74 \pm 1.42$	$48.47 \pm 1.98$
Coreset	$26.58 \pm 0.29$	$\textbf{31.31} \pm \textbf{1.48}$	$35.46 \pm 0.75$	$40.05\pm1.51$	$43.86 \pm 1.33$	$46.15\pm1.21$	$49.38 \pm 0.70$
VAAL	$26.58 \pm 0.29$	$29.37 \pm 1.73$	$31.61 \pm 2.18$	$39.66 \pm 1.88$	$42.49 \pm 2.04$	$46.01 \pm 1.54$	$\textbf{49.73} \pm \textbf{0.69}$
DBAL	$26.58 \pm 0.29$	$30.17 \pm 1.18$	$\textbf{36.37} \pm \textbf{0.99}$	$39.77 \pm 1.83$	$44.22 \pm 0.76$	$45.79 \pm 0.40$	$48.28 \pm 0.66$
QBC	$26.58 \pm 0.29$	$29.33 \pm 0.50$	$34.35 \pm 2.58$	$\textbf{40.55} \pm \textbf{1.34}$	$43.06 \pm 2.16$	$45.37 \pm 0.77$	$48.09 \pm 0.84$
BALD	$26.58 \pm 0.29$	$30.02\pm0.98$	$35.04\pm1.25$	$40.34\pm1.32$	$43.65\pm1.05$	$46.14 \pm 1.73$	$49.10\pm1.28$

Table 15. Mean Accuracy and Standard Deviation on CIFAR10/100 test set with annotation size as 5% of training set. Results reported are averaged over 5 runs.