Noise-free Loss Gradients: A Surprisingly Effective Baseline for Coreset Selection

Saumyaranjan Mohanty

Department of Artificial Intelligence Indian Institute of Technology Hyderabad

Chimata Anudeep Department of Computer Science BITS Pilani KK Birla Goa Campus

Konda Reddy Mopuri

Department of Artificial Intelligence Indian Institute of Technology Hyderabad

Reviewed on OpenReview: https://openreview.net/forum?id=OE4P1tW8iQ

Abstract

The exponential rise in size and complexity of deep learning models and datasets have resulted in a considerable demand for computational resources. Coreset selection is one of the methods to alleviate this rising demand. The goal is to select a subset from a large dataset to train a model that performs almost at par with the one trained on the large dataset while reducing computational time and resource requirements. Existing approaches either attempt to identify remarkable samples (e.g., Forgetting, Adversarial Deepfool, EL2N, etc.) that stand out from the rest or solve complex optimization (e.g., submodular maximization, OMP) problems to compose the coresets. This paper proposes a novel and intuitive approach to efficiently select a coreset based on the similarity of loss gradients. Our method works on the hypothesis that gradients of samples belonging to a given class will point in similar directions during the early training phase. Samples with most neighbours that produce similar gradient directions, in other words, that produce noise-free gradients, will represent that class. Through extensive experimentation, we have demonstrated the effectiveness of our approach in out-performing state-of-the-art coreset selection algorithms on a range of benchmark datasets from CIFAR-10 to ImageNet with architectures of varied complexity (ResNet-18, ResNet-50, VGG-16, ViT). We have also demonstrated the effectiveness of our approach in Generative Modelling by implementing coreset selection to reduce training time for various GAN models (DCGAN, MSGAN, SAGAN, SNGAN) for different datasets (CIFAR-10, CIFAR-100, Tiny ImageNet) while not impacting the performance metrics significantly. Source code is provided at URL.

1 Introduction

Large-scale datasets have become essential for training state-of-the-art deep learning models. All the application fronts of Artificial Intelligence, such as Computer Vision, Natural Language, and Speech Processing, have accumulated massive datasets. Concurrently, the complexity of deep learning models has grown steadily. Modern deep learning tasks (especially supervised ones) require extensive hyperparameter tuning during training to achieve the best performance. These have resulted in an exponential rise in computational requirement (Schwartz et al., 2019) and carbon footprint(Heikkilä; Killamsetty et al., 2021a).

ai 23 resch04001 @iith.ac. in

chanudeep03@gmail.com

krmopuri@ai.iith.ac.in

Coreset selection aims to mitigate the above issues by finding the most representative data samples from the given training set. Mainly, coreset selection attempts to approximate the learning characteristics of the complete data (e.g., the loss function) (Feldman, 2020). By obtaining a representative coreset with a cardinality of a fraction of the entire dataset, training duration and computational requirement for end-to-end training can be reduced significantly while delivering the desired generalization performance.

Forgetting(Toneva et al., 2019) attempts to identify data points that are difficult to learn as the coreset. The idea presented is that a model trained on "hard to classify" samples naturally generalizes onto the simpler samples. Similarly, training loss based methods attempt to select samples that contribute more to the training of deep neural networks. Gradient Norm (GraNd) and the Error L2-Norm (EL2N) scores introduced by Paul et al. (2023) prune significant fractions of training data without sacrificing much of the classifier's generalization performance. Section 2 comprehensively summarizes the existing works.

Multiple lines of thought based on loss gradients have emerged to compose coresets. Gradient matching based methods aim to compose a coreset whose gradients (weighted combination) closely approximate the gradients produced by the entire training dataset. For instance, CRAIG (Mirzasoleiman et al., 2020) suggests finding the optimal coreset by converting the gradient matching problem to a submodular function maximization and using a greedy approach to optimize it. Because of the individual weights (or learning rates) on the samples in CRAIG, it may not be straightforward to train the models on the coresets in the mini-batch gradient descent framework. Hence, the training speed gains are not maximal on such coresets.

Similarly, Gradmatch(Killamsetty et al., 2021a) constructs an objective for matching the gradients computed over the coreset with that calculated over the complete dataset. They propose minimizing the matching error by casting the objective as a weakly submodular maximization problem and solving it using an orthogonal matching pursuit (OMP) based greedy algorithm. As it identifies the coreset by selecting different subsets during model training to perform gradient matching, training must be carried out each time to obtain coresets at various fractions of the entire dataset.

In this paper, we propose a novel methodology for selecting a coreset based on the hypothesis that most samples from the same class produce gradients in similar directions (particularly in the early stages of the model training, elaborated further in Section 3, Figures 1 and 2). Our method proposes that samples with a substantial gradient similarity to many others belonging to the same class would be ideal candidates for constituting the representative coreset. This approach differs from the existing coreset ideas driven by 'distinctness' or 'difficulty', which prefer samples that give rise to noisy (or minority) loss gradients. Hence, we refer to the proposed method as 'Noise-free Gradients' for coreset selection. In the case of a multi-class classification task, we identify representative samples class-wise and combine them into a coreset.

Our approach uses cosine similarity between gradient vectors to rank data samples according to their ability to represent other samples. Based on the desired coreset size, the top-ranked samples from each class are picked to represent the whole dataset. This makes the proposed approach computationally attractive compared to the optimization-driven coreset ideas (Refer to Section 2 for a detailed review of such works). Moreover, while the existing approaches, such as CRAIG, use a per-sample learning rate while training from the coreset, our method uses single learning rate for all the samples. We also introduce various steps to reduce the computational time. We have applied this coreset selection methodology to two popular and diverse application areas, namely, 1) Image classification and 2) Generative Modelling. The details are discussed in Section 4.

In summary, the major contributions of our work can be summarized as:

- Contrary to the existing 'difficulty/distinctness' notions, this paper introduces a novel, intuitive coreset selection method for identifying class-wise representative samples driven by noise-free loss gradients.
- Also, compared to the existing gradient-based coreset selection methods that are computationally demanding (e.g., submodular function optimization), our method is computationally efficient.
- We thoroughly evaluate the proposed method over multiple object recognition datasets of varying complexity. Experiments with different deep neural network classifiers (CNN and Vision Trans-

former) and cross-architecture generalization studies demonstrate that our method achieves higher accuracy than state-of-the-art coreset selection methods.

• We apply the proposed method in the context of GAN training and demonstrate consistent performance improvement over existing methods.

2 Related Works

Multiple coreset selection methods have recently been proposed. Here, we briefly present the most prominent works.

Toneva et al. (2019) proposed coreset selection through catastrophic forgetting. During the training, they defined a flip in classifying a training example from the correct label to an incorrect one as a 'forgetting' event. Removing unforgettable examples results in minimal performance drop while speeding up the training and hyperparameter tuning.

K-Center Greedy approximation (Sener & Savarese, 2018) attempts to solve the minimax facility location problem to select coresets from a large dataset such that the maximum distance between points in the non-coreset and its closest point in the coreset is minimized.

Uncertainty-based methods work on the idea that samples having lower confidence may have a higher impact during training than those with higher confidence. Thus, these methods suggest constituting the coresets with the samples with lower confidence. Commonly used metrics to calculate sample uncertainty are least confidence (Shen et al., 2018), entropy (Settles, 2012), and margin (Coleman et al., 2020).

Adversarial Deepfool (Ducoffe & Precioso, 2018) and Contrastive Active Learning (Margatina et al., 2021) work to find data points distributed near the decision boundary. Samples that lie close to the classification boundary are treated as difficult samples to learn. Error-based methods try to select samples that contribute more to the training of the neural networks. Two metrics called the Gradient Normed (GraNd) and the Error L2-Norm (EL2N) scores are introduced by Paul et al. (2023) that help in pruning significant fractions of training data without sacrificing test accuracy. *GraNd* measures the importance of each sample to the training loss at early epochs. *EL2N* approximates the *GraNd* score, which measures the norm of the error vector, with a higher score indicating higher potential influence. The authors conclude that the images with higher scores tend to be harder to learn (forgettable examples). Their method chooses these forgettable samples as the coreset.

RETRIEVE (Killamsetty et al., 2021d) formulates coreset selection for Semi-Supervised Learning (SSL) as a bi-level optimization problem. This method considers labelled and unlabeled sets to develop the bi-level optimization problem. It uses a greedy algorithm to select the coreset that minimizes the labelled set loss. GLISTER (Killamsetty et al., 2021b), Generalization-based Data Subset Selection for Efficient and Robust Learning, applies bi-level optimization for supervised and active learning. It formulates coreset selection as a bi-level optimization problem that maximizes the log-likelihood on a held-out validation dataset.

Gradient matching-based methods work on the expectation that the optimal coreset can approximate the gradients produced by the entire training dataset. CRAIG (Mirzasoleiman et al., 2020) selects representative subsets that closely approximate the entire gradient. They achieve this by converting the gradient matching problem to optimizing a submodular function using a greedy approach. *Gradmatch* (Killamsetty et al., 2021a) method follows a similar approach. It introduces a squared L_2 regularization term over the weight vectors and uses a greedy Orthogonal Matching Pursuit (OMP) algorithm to select the coreset iteratively. After training on the resulting coreset for a pre-determined number of epochs, the algorithm repeats the coreset construction using the latest iterate.

Several methods have utilized submodular functions (Liu, 2020) to leverage their ability to measure diversity and information, then use greedy algorithms to maximize a submodular function for coreset selection. Iyer et al. (2021) have studied various submodular functions such as Graph Cut, entropy, and facility location and utilize greedy optimization algorithms for data summarization. Xia et al. (2023) have proposed a universal method termed 'Moderate Coreset' to tackle the task-specific nature of dataset selection. They calculate the distance between the hidden representation of samples and the representational class centres. Based on these Euclidean distances, they rank the data points in an ascending order and select the data points closest to the distance median as a coreset.

Yang et al. (2023) proposed a parameter influence-based optimization for dataset selection. They argue that not all training data contributes uniformly to the model's parameter learning. So, an influence function-based iterative method can be designed to linearly estimate the impact of omitting a specific subset on the model parameters. However, they have not considered lower coreset sizes than 50% of the original dataset size.

Sinha et al. (2019) introduced the application of coreset selection to sub-sample a larger batch to produce a smaller batch to speed up GAN training significantly while not sacrificing performance to a significant level. DeVries et al. (2020) introduced instance selection based upon the high-density region of the data manifold to speed up GAN training.

Guo et al. (2022) has developed an excellent and comprehensive code library named *DeepCore* that implements current popular and state-of-the-art coreset selection methods in a unified framework based on PyTorch (Paszke et al., 2019). DeepCore framework has provided a unified comparative analysis among 15 algorithms for CIFAR-10 and 14 algorithms for ImageNet-1K. The algorithms compared are: Contextual Diversity (Agarwal et al., 2020), Herding (Welling, 2012), k-Center Greedy (Sener & Savarese, 2018), Least Confidence, Entropy, Margin (Coleman et al., 2020), Forgetting (Toneva et al., 2019), GraNd (Paul et al., 2023), Cal (Margatina et al., 2021), DeepFool (Ducoffe & Precioso, 2018), Craig (Mirzasoleiman et al., 2020), GradMatch (Killamsetty et al., 2021b), Glister (Killamsetty et al., 2021b), Facility Location and GraphCut (Iyer et al., 2021). Due to its large running time, the DeepFool algorithm is not considered for the ImageNet-1K dataset. We have also utilized publicly available implementation of Moderate Coreset (Xia et al., 2023) for CIFAR-10, CIFAR-100 and Tiny ImageNet datasets and presented the comparative analysis in Section 5.

3 Methodology

This section provides a detailed description of the proposed *Noise-free Gradients* approach. We first focus on the supervised classification task of object recognition using Deep Neural Network classifiers (CNNs and Vision Transformers) because of their proven effectiveness. We have also applied our proposed method to speed up GAN training without sacrificing the performance of the model significantly. Table 1 provides the notation used throughout the paper.

Symbol	Description
$\mathbb{V} = \{(x_i, y_i)\}$	Large Training Dataset
$\mathbb{S}\subset\mathbb{V}$	Coreset of \mathbb{V} (target)
θ	Parameters of the classifier
Φ	Threshold on the gradient similarity for neighborhood identification
\mathbb{V}_{c}	Training data belonging to class c
$g^{ heta}_{x_i}$	Loss gradients computed for data sample (x_i, y_i) at the last fully connected layer
$\ x\ $	l_2 norm of x
$ x _1$	l_1 norm of x
$ \mathbb{A} $	Cardinality of set \mathbb{A}
1	Indicator function
< .,. >	Dot product operator

Table	1:	Notation

Our objective is to select a representative subset S of the complete dataset V such that model θ^{S} trained on S has a generalization performance close to that of the model θ^{V} trained on V. Training deep neural networks

is reduced to an empirical risk minimization problem often optimized in the gradient descent framework. In practice, the incremental Gradient (IG) methods, such as Stochastic Gradient Descent (SGD), iteratively estimate the gradient on mini-batches of training data that construct the parameter updates.

Existing gradient-based coreset selection methods such as CRAIG(Mirzasoleiman et al., 2020) and Gradmatch(Killamsetty et al., 2021a) try to find an optimal coreset such that the weighted sum of the gradients of the coreset elements remains within an error margin of the gradients of the entire dataset. The objective function can be written as:

$$\underset{w,\mathbb{S}}{\operatorname{arg\,min}} F(\frac{1}{|\mathbb{V}|} \sum_{(x_i,y_i)\in \mathbb{V}} g_{x_i}^{\theta}, \frac{1}{||w||} \sum_{(x_i,y_i)\in \mathbb{S}} w_{x_i} g_{x_i}^{\theta})$$
(1)

Where w is the vector of weights associated with the elements of the coreset, and F is a distance metric. It can be re-written as:

$$\mathbb{E}_{\theta}\left[\frac{1}{|\mathbb{V}|}\sum_{(x_i,y_i)\in\mathbb{V}}g_{x_i}^{\theta}\right] = \mathbb{E}_{\theta}\left[\frac{1}{||w||}\sum_{(x_i,y_i)\in\mathbb{S}}w_{x_i}g_{x_i}^{\theta}\right] + \epsilon$$
(2)

Where ϵ is the error term having the same dimension as the gradient vector, note that equation 2 considers the expected value of the approximation error in the parameter space.

Our method identifies the representative samples based on their gradients. Intuitively, data samples whose gradients are similar to most other samples best approximate the complete dataset. In other words, these representative samples result in local minima close to the minima achieved by the entire dataset.

Our method thus selects a subset S of a desired cardinality that closely approximates the whole dataset V. We define the normalized gradient similarity between two samples as

$$\rho(x_i, y_i, x_j, y_j, \theta) = \frac{\langle g_{x_i}^{\theta}, g_{x_j}^{\theta} \rangle}{\|g_{x_i}^{\theta}\| \|g_{x_j}^{\theta}\|}$$
(3)

For each sample x_i in the dataset \mathbb{V} , we can measure its ability to represent \mathbb{V} as

$$f(x_i) = \mathbb{E}_{\theta} \left[\sum_{x_j \in \mathbb{V}, \ j \neq i} \rho(x_i, y_i, x_j, y_j, \theta) \right]$$
(4)

Using this measure, we compute the suitability of every sample in \mathbb{V} to become an element of the coreset \mathbb{S} . Essentially, this translates to sorting the dataset samples in the decreasing order of this measure and composing the coreset of a desired cardinality. The exact steps involved in this process are described in the next paragraph.

We start with a randomly initialized model and update its parameters on the complete dataset for a few epochs. In our experiments, we observe that generally, 5 to 10 epochs (during which the loss value doesn't plateau) are sufficient. We save the checkpoints of the model parameters after each epoch. For each of these checkpoints, we calculate the gradients of the loss function with respect to the model parameters computed at each data sample. We score the dataset samples based on their ability to represent other samples as denoted in equation 4. We aggregate these scores across the saved checkpoints as an approximation to the expectation over the parameter space.

Figure 1 shows the histogram of pairwise cosine similarity of gradients of images belonging to airplane class for CIFAR-10 dataset on ResNet-18 architecture for initial four epochs. As we can see, during initial training, the most images have gradients with high cosine similarity values, indicating moving in similar directions. Figure 2 shows the histogram of pairwise cosine similarity of gradients of images belonging to (i) airplane



Figure 1: Histogram of pairwise cosine similarity of gradients during the training of ResNet-18 for the initial four epochs. Samples belong to the 'airplane' class of CIFAR-10 dataset. A strong similarity (or alignment) of the gradients can be observed.



Figure 2: histogram of pairwise cosine similarity of gradients of images belonging to (i) airplane class (top panel), (ii) airplane and bird class (middle panel), and (iii) airplane and truck class (bottom panel). While intra-class samples display strong similarity, inter-class samples have gradients in different directions.

class (top panel), (ii) airplane and bird class (middle panel), and (iii) airplane and truck class (bottom panel). All the computations are done for the model saved after the first epoch. It indicates that images belonging to the same class will have gradients in similar directions, while images belonging to different classes will have gradients in different directions. This behaviour forms the basis of our method for computation of the coreset.

To improve the time efficiency of our approach, we slightly modify it to utilize a nearest-neighbor search algorithm. We utilize radius-based neighbour learning to obtain the local density of data points. The persample scores denoted by equation 4 are measured in terms of the number of dataset samples present nearby in the gradient space. A nearest neighbour search algorithm finds the number of samples within a given radius from each sample (in other words, with similarity more than a threshold Φ) and assigns it as its score. We then rank the samples according to their aggregated scores across multiple checkpoints. The top-ranked images are selected as the representative coreset.

The formulation of *Noise-free Gradients* can be represented as shown in equation 5. For a given class c, the top-ranked samples x_j are selected based on their scores.

$$x_j = \operatorname*{arg\,max}_{(x_i, y_i) \in \mathbb{V}_c} \sum_{\theta} \sum_{j \neq i} \mathbb{1}(\rho(x_i, y_i, x_j, y_j, \theta) > \Phi)$$
(5)

Algorithm 1 presents our approach more formally.

Algorithm 1 Noise-free Gradients for Coreset algorithm

Require: Train set: \mathbb{V} ; Total epochs: T; number of classes: C; number of coreset images per class: N; Model checkpoint after initial T epochs : θ

Ensure: Coreset \mathbb{S}

```
1: for class c in 1,..., C do
            for (x_i, y_i) \in \mathbb{V}_c do
2:
                  for epochs t in 1,..., T do
 3:
                        compute g_{x_i}^{\theta_t}
 4:
                 end for

f(x_i) = \sum_{\theta_t} \sum_{(x_j, y_j) \in \mathbb{V}_c, j \neq i} \mathbb{1}(\rho(x_i, y_i, x_j, y_j, \theta_t) > \Phi)
 5:
 6:
                  Store f(x_i)
 7:
 8:
            end for
9: end for
10: \mathbb{S} = \emptyset
11: for class c in 1,..., C do
            \mathbb{S} \leftarrow \mathbb{S} \cup \operatorname{argsort}_{(x_i, y_i) \in \mathbb{V}_c} f(x_i) [: N]
12:
13: end for
```

▷ Descending order

4 Implementation of Noise-free Gradients

We have implemented the proposed method *Noise-free Gradients* using the PyTorch framework.

4.1 Gradients with respect to the Classification layer

Deep learning models such as ResNet (He et al., 2015; 2016), VGG (Simonyan & Zisserman, 2015), and ViT Small (Lee et al., 2021) have millions of parameters, and it is impractical to consider the gradient of loss function with respect to each of the parameters. It is observed that the gradient mostly captures the variation of gradient norm with respect to the parameters of the last (classification) layer of the neural network (Katharopoulos & Fleuret, 2019). Similar to the earlier works (Ash et al., 2020; Killamsetty et al., 2021a), we avoid computation of the gradient of the loss function with respect to all the model parameters, restricting to only the final fully connected layer while calculating the score of data samples (equation 4).

4.2 Nearest Neighbor algorithm

The complexity of similarity score computation (equation 4) among the samples of a particular class of size N is $\mathcal{O}(N^2)$. The complexity and computation time are linear in the number of classes in the dataset C. Instead of computing N similarity scores for each sample in a given class, our approach finds the number of these N samples within a close neighborhood. This vital modification saves nontrivial complexity. We have utilized the radius-based nearest neighbor algorithm implementation from scikit-learn (Pedregosa et al.,

2011) library with the proposed distance (or similarity) measure mentioned in equation 3 as the metric for identification of neighbors within a given radius. Pairwise cosine similarity function of torchmetrics library¹ is used. This results in a $\approx 30X$ speed up of the computation time compared to the naive method. We have used a threshold value (Φ) of 0.2 in all our experiments. For the ablation experiments on the effect of this threshold, readers can refer to the supplementary document (Section 2).

5 Experiments and Results

5.1 Experimental Setup

Applications. We have evaluated our coreset selection methodology on image classification and GAN training.

Datasets. For image classification application, we evaluate the effectiveness of our method on four popularly used benchmark object recognition datasets, i.e., CIFAR-10(Krizhevsky, 2009), CIFAR-100(Krizhevsky, 2009), Tiny ImageNet (Le & Yang, 2015) and ImageNet-1K (Russakovsky et al., 2015). CIFAR-10 dataset consists of 50,000 colour images of dimension $32 \times 32 \times 3$ from 10 different classes, each class having 5,000 images. CIFAR-100 dataset consists of 50,000 training images from 100 classes with 500 training images per class. Tiny ImageNet dataset consists of 100,000 training images from 200 classes with 500 training images per class. ImageNet-1K is a subset of the larger dataset ImageNet, an image dataset organized according to the WordNet hierarchy. ImageNet-1K consists of 1000 classes, with 1, 281, 167 training images and 50,000 validation images. In Summary, to test the robustness of our method, we have considered four popular classification datasets with the number of classes varying from 10 to 1000.

Classifier. We have used randomly initialized ResNet-18 (He et al., 2016), ResNet-50 (He et al., 2015), VGG-16 (Simonyan & Zisserman, 2015) and ViT Small (Lee et al., 2021) architectures to carry out comparative performance analysis of existing SOTA algorithms against our method. The first three are CNN architectures that vary in complexity and number of parameters, and the fourth is a vision transformer. We have also carried out cross-architecture generalization to demonstrate the effectiveness of our method.

GAN models. For GAN training application, we evaluate the effectiveness of our method on four different GAN architectures i.e. DCGAN (Goodfellow et al., 2014), MSGAN (Mao et al., 2019), SNGAN Miyato et al. (2018) and SAGAN Zhang et al. (2019).

GAN Metrics. We have considered the Inception Score (IS) (Salimans et al., 2016) and Fretchet Inception Distance (FID) (Heusel et al., 2018) for comparing the performance of our method against existing coreset selection for GAN methodologies.

Baselines We have considered the results reported by DeepCore benchmark wherever available and utilized DeepCore implementation to carry out experiments wherever results are not available.² We have utilized publicly available implementation of Moderate-DS³ and instance selection for GAN⁴. Abbreviations used in results for different methods which are part of DeepCore are: GC(GraphCut), F(Forgetting), R(Random), FL(Facility Location), GN(GraNd), CR(CRAIG), H(Herding).

Implementation Experiments for CIFAR-10 and CIFAR-100 are carried out for ten individual runs with different random seeds. Experiments for tiny ImageNet were carried out for five separate runs, and ImageNet1K were carried out for two individual runs due to substantial computational requirements. In each iteration, number of training epochs were set at 200, and an SGD optimizer with a momentum of 0.9 was used with a batch size of 128 and a learning rate of 0.1. For ImageNet-1K, batch size of 256 was used.

For GAN training, we have used default settings for all the GAN models recommended by their respective authors. For instance selection, a 50% retention ratio is used. For our method, a 50% coreset selection percentage is used. FID scores are lower the better, and Inception scores are higher the better. The number of steps used for full dataset training is 200k, and the steps used for coreset-based training is 100k.

¹https://lightning.ai/docs/torchmetrics/stable/pairwise/cosine_similarity.html

²https://github.com/PatrickZH/DeepCore

³https://github.com/tmllab/2023_ICLR_Moderate-DS

 $^{{}^{4} \}tt{https://github.com/uoguelph-mlrg/instance_selection_for_gans$

5.2 Results for Image Classification

5.2.1 Results for CIFAR-10

A comparison of accuracy values obtained from DeepCore, Moderate-DS, and *Noise-free Gradients* on the CIFAR-10 dataset for various fractions of the original dataset is presented in Table 2, Table 3, Table 4 and Table 5 for ResNet-18, ResNet-50, VGG-16 and ViT Small architecture respectively. The algorithm with the best accuracy at each percentage level reported by DeepCore is mentioned within brackets against the accuracy value. The average rank of a method is calculated as the mean of its rank across different fractions of the dataset considered. As can be observed from the values, our method outperforms all the existing SOTA algorithms across architectures of diverse complexity and depth.

Table 2: Comparison of results for CIFAR-10 on ResNet-18

Size	DeepCore	Moderate-DS	Our Method
$\begin{array}{c} 0.5\% \\ 1.0\% \\ 5.0\% \\ 10.0\% \\ 20.0\% \\ 30.0\% \end{array}$	$ \begin{array}{c} 34.90 \pm 2.30 \; ({\rm GC}) \\ 42.80 \pm 1.30 \; ({\rm GC}) \\ 65.70 \pm 1.20 \; ({\rm GC}) \\ 76.60 \pm 1.50 \; ({\rm GC}) \\ 87.10 \pm 0.50 \; ({\rm R}) \\ 91.70 \pm 0.3 \; ({\rm F}) \end{array} $	$\begin{array}{c} 33.93 \pm 1.10 \\ 39.83 \pm 0.56 \\ 69.20 \pm 1.30 \\ 79.10 \pm 0.12 \\ 86.49 \pm 0.29 \\ 89.39 \pm 0.12 \end{array}$	$\begin{array}{c} 43.64 \pm 0.87 \\ 54.06 \pm 0.92 \\ 72.96 \pm 0.28 \\ 79.12 \pm 0.40 \\ 87.18 \pm 0.45 \\ 89.62 \pm 0.50 \end{array}$
Rank	2.17	2.67	1.17

Table 3: Comparison of results for CIFAR-10 on ResNet-50

Size	DeepCore	Moderate-DS	Our Method
0.5% 1.0% 5.0% 10.0%	$\begin{vmatrix} 32.85 \pm 1.20 & (GC) \\ 36.58 \pm 1.50 & (GC) \\ 63.31 \pm 1.53 & (GC) \\ 70.64 \pm 1.10 & (E) \end{vmatrix}$	13.73 ± 4.82 18.86 ± 1.55 53.89 ± 1.03 65.47 ± 1.82	$47.32 \pm 0.69 \\ 56.06 \pm 0.10 \\ 73.04 \pm 0.32 \\ 78.06 \pm 0.41$
10.0% 20.0% 30.0%	$\begin{array}{c} 70.64 \pm 1.10 \text{ (F)} \\ 85.06 \pm 1.48 \text{ (FL)} \\ \textbf{89.83} \pm \textbf{0.30 (F)} \end{array}$	$\begin{array}{c} 05.47 \pm 1.83 \\ 84.00 \pm 0.79 \\ 88.74 \pm 0.27 \end{array}$	$\begin{array}{c} 78.06 \pm 0.41 \\ 86.11 \pm 0.48 \\ 89.22 \pm 0.20 \end{array}$
Rank	1.83	3.00	1.17

Table 4: Comparison of results for CIFAR-10 on VGG-16

Size	DeepCore	Moderate-DS	Our method
$\begin{array}{c} 0.5\% \\ 1.0\% \\ 5.0\% \\ 10.0\% \\ 20.0\% \\ 30.0\% \end{array}$	$\begin{array}{c} 32.85 \pm 1.2 \ ({\rm GC}) \\ 36.58 \pm 1.5 \ ({\rm GC}) \\ 63.31 \pm 1.53 \ ({\rm GC}) \\ 70.64 \pm 1.1 \ ({\rm F}) \\ 85.06 \pm 1.48 \ ({\rm FL}) \\ \mathbf{89.83 \pm 0.3 \ (F)} \end{array}$	$\begin{array}{c} 13.73 \pm \!$	$\begin{array}{c} 47.32 \pm 0.69 \\ 56.06 \pm 0.10 \\ 73.04 \pm 0.32 \\ 78.06 \pm 0.41 \\ 86.11 \pm 0.48 \\ 89.22 \pm 0.2 \end{array}$
Rank	1.67	3	1.33

Table 5:	Comparison	of results	\mathbf{for}	CIFAR-10 or	ı ViT small
----------	------------	------------	----------------	-------------	-------------

Size	DeepCore	Moderate-DS	Our Method
0.5% 1.0% 5.0%	$\begin{array}{c} 30.05 \pm 1.21 \; (\text{GC}) \\ 37.31 \pm 0.52 \; (\text{GC}) \\ 53.96 \pm 0.45 \; (\text{GC}) \end{array}$	$\begin{array}{c} 23.10 \pm 1.41 \\ 27.00 \pm 4.80 \\ 51.30 \pm 2.30 \end{array}$	$35.22 \pm 0.31 \\ 38.50 \pm 0.27 \\ 56.82 \pm 0.20$
10.0% 20.0% 30.0% 40.0%	57.43 ± 0.69 (GC) 64.98 ± 0.90 (GC) 67.38 ± 1.01 (GC) 71.24 ± 0.20 (GC)	54.40 ± 2.35 66.10 ± 1.66 70.40 ± 0.28 72.60 ± 1.40	58.32 ± 0.16 67.16 ± 0.52 71.43 ± 0.49 75.45 ± 0.70
40.0% 50.0% 60.0%	$71.24 \pm 0.20 \text{ (GC)}$ $72.98 \pm 2.10 \text{ (F)}$ $77.68 \pm 0.20 \text{ (GC)}$	$72.00 \pm 1.40 \\ 73.10 \pm 4.09 \\ 76.95 \pm 1.85$	73.43 ± 0.10 78.07 ± 0.68 77.11 ± 0.10
Rank	2.33	2.55	1.11

5.2.2 Results for CIFAR-100

A comparison of accuracy values obtained from DeepCore, Moderate-DS, and *Noise-free Gradients* on the CIFAR-100 dataset for various fractions of the original dataset is presented in Table 6, Table 7, Table 8 and Table 9 for ResNet-18, ResNet-50, VGG-16 and ViT Small architectures respectively. The algorithm with the best accuracy at each percentage level reported by DeepCore is mentioned within brackets against the accuracy value. The average rank of a method is calculated as the mean of its rank across different fractions of the dataset considered. We outperform all the other methods for training on ResNet-50, VGG-16 and ViT Small architecture while performing better for lower selection percentages (up to 20%) for training on ResNet-18 architecture.

5.2.3 Results for Tiny ImageNet and ImageNet-1K

We believe the effectiveness of coreset selection algorithms must be tested in the face of complex datasets. Hence, we evaluate our algorithm with the Tiny ImageNet and ImageNet-1K datasets. Similar to other experiments, we train multiple models on the resulting coresets. Each model is trained for 200 epochs with a random PyTorch seed. The comparative analysis is tabulated in Table 10, 11, 12 and 13. The algorithms with the best accuracy at each percentage level reported by DeepCore are mentioned within brackets against the accuracy value. The average rank of a method is calculated as the mean of its rank across different fractions of the dataset considered. Our method consistently outperforms all the other coreset selection algorithms. Table 6: Comparison of results for CIFAR-100 on ResNet-18

Size	DeepCore	Moderate-DS	Our Method
$\begin{array}{c} 0.5\% \\ 1.0\% \\ 5.0\% \\ 10.0\% \\ 20.0\% \\ 30.0\% \\ 40.0\% \end{array}$	$\begin{array}{c} 8.49 \pm 0.68 \ ({\rm GC}) \\ 13.04 \pm 0.22 \ ({\rm GC}) \\ 31.18 \pm 0.37 \ ({\rm GC}) \\ 41.98 \pm 0.88 \ ({\rm GC}) \\ 56.68 \pm 0.98 \ ({\rm GC}) \\ 64.05 \pm 0.28 \ ({\rm GC}) \\ 69.37 \pm 0.30 \ ({\rm GC}) \end{array}$	$\begin{array}{c} 5.36 \pm 0.12 \\ 7.78 \pm 0.20 \\ 18.62 \pm 0.33 \\ 32.57 \pm 1.40 \\ 54.4 \pm 0.19 \\ 62.97 \pm 0.41 \\ 67.33 \pm 0.22 \end{array}$	$\begin{array}{c} 12.0 \pm 0.11 \\ 18.95 \pm 0.08 \\ 37.67 \pm 0.10 \\ 46.48 \pm 0.28 \\ 58.06 \pm 0.34 \\ 62.58 \pm 0.24 \\ 66.31 \pm 0.16 \end{array}$
60.0%	$ $ 73.80 \pm 0.31 (GN)	72.02 ± 0.30	70.74 ± 0.20
Rank	1.63	2.63	1.75

Table 8: Comparison of results for CIFAR-100 on VGG-16

Size	DeepCore	Moderate-DS	Our Method
$\begin{array}{c} 0.5\% \\ 1.0\% \\ 5.0\% \\ 10.0\% \\ 20.0\% \\ 30.0\% \\ 40.0\% \\ 50.0\% \end{array}$	$\begin{array}{c} 2.45 \pm 0.57 \ ({\rm GC}) \\ 4.47 \pm 1.28 \ ({\rm GC}) \\ 27.23 \pm 0.70 \ ({\rm GC}) \\ 38.17 \pm 1.32 \ ({\rm GC}) \\ 52.23 \pm 0.47 \ ({\rm GC}) \\ 58.87 \pm 0.45 \ ({\rm GC}) \\ 62.52 \pm 0.23 \ ({\rm GC}) \\ 62.64 \pm 0.21 \ ({\rm GC}) \end{array}$	$\begin{array}{c} 2.00 \pm 0.10 \\ 2.10 \pm 0.23 \\ 13.70 \pm 0.91 \\ 28.10 \pm 1.60 \\ 49.90 \pm 0.30 \\ 59.50 \pm 0.30 \\ 64.92 \pm 0.93 \\ 64.91 \pm 0.60 \end{array}$	$\begin{array}{c} 8.92 \pm 0.45 \\ 15.2 \pm 0.46 \\ 32.12 \pm 0.58 \\ 44.67 \pm 0.05 \\ 53.95 \pm 0.30 \\ 61.23 \pm 0.17 \\ 66.17 \pm 0.23 \\ 0.23 \end{array}$
$50.0\%\ 60.0\%$	$68.62 \pm 1.21 (GC)$ 73.34 \pm 0.39 (F)	69.11 ± 0.62 71.87 ± 0.91	$ \begin{array}{r} 69.47 \pm 0.19 \\ 71.82 \pm 0.20 \end{array} $
Rank	2.22	2.55	1.22

Table 10: Comparison of results for Tiny ImageNet on ResNet-18

Size	DeepCore	Moderate-DS	Our Method
$\begin{array}{c} 0.5\% \\ 1.0\% \\ 5.0\% \\ 10.0\% \\ 20.0\% \\ 30.0\% \\ 40.0\% \end{array}$	$\begin{array}{c} 3.29 \pm 0.23 \ ({\rm FL}) \\ 5.34 \pm 0.45 \ ({\rm GC}) \\ 13.73 \pm 0.71 \ ({\rm GC}) \\ 21.23 \pm 0.16 \ ({\rm GC}) \\ 29.97 \pm 0.72 \ ({\rm GC}) \\ 33.73 \pm 0.37 \ ({\rm GC}) \\ 35.74 \pm 0.57 \ ({\rm FL}) \end{array}$	$\begin{array}{c} 3.34 \pm 0.33 \\ 4.11 \pm 0.24 \\ 14.03 \pm 1.01 \\ 21.40 \pm 0.41 \\ 30.99 \pm 0.93 \\ 37.18 \pm 0.87 \\ 40.21 \pm 0.24 \end{array}$	$\begin{array}{c} 6.91 \pm 0.24 \\ 11.09 \pm 0.45 \\ 22.83 \pm 0.98 \\ 28.05 \pm 0.82 \\ 34.91 \pm 0.43 \\ 40.36 \pm 0.43 \\ 40.41 \pm 0.17 \end{array}$
60.0% Rank	$\frac{39.74 \pm 0.73 \text{ (F)}}{2.75}$	$\frac{46.19 \pm 0.64}{2.12}$	44.74 ±0.44 1.12

Table 12: Comparison of results for Tiny ImageNet on ViT small

Size	DeepCore	Moderate-DS	Our Method
$\begin{array}{c} 0.5\% \\ 1.0\% \\ 5.0\% \\ 10.0\% \\ 20.0\% \\ 30.0\% \\ 40.0\% \\ 50.0\% \end{array}$	$\begin{vmatrix} 2.88 \pm 0.04 \text{ (FL)} \\ 4.23 \pm 0.25 \text{ (F)} \\ 8.41 \pm 0.15 \text{ (GC)} \\ 11.53 \pm 1.23 \text{ (GC)} \\ 12.54 \pm 0.06 \text{ (GC)} \\ 14.09 \pm 0.11 \text{ (GC)} \\ 15.69 \pm 0.45 \text{ (GC)} \\ 17.34 \pm 1.05 \text{ (GC)} \end{vmatrix}$	$\begin{array}{c} 3.35 \pm 0.20 \\ 4.34 \pm 0.13 \\ 13.56 \pm 0.23 \\ 15.84 \pm 0.22 \\ 22.03 \pm 0.32 \\ 25.35 \pm 0.31 \\ 28.40 \pm 0.24 \\ 30.68 \pm 0.14 \end{array}$	$5.92 \pm 0.12 \\9.18 \pm 0.25 \\18.10 \pm 0.15 \\22.58 \pm 0.06 \\28.26 \pm 0.31 \\33.66 \pm 0.21 \\37.29 \pm 0.32 \\40.92 \pm 0.32 \\$
60.0%	$17.34 \pm 1.05 (GC)$ $19.23 \pm 0.32 (GC)$	32.74 ± 0.27	40.32 ± 0.30 43.80 ± 0.38
Rank	3.00	2.00	1.00

Table 7: Comparison of results for CIFAR-100 on ResNet-50

Size	DeepCore	Moderate-DS	Our Method
0.5%	$7.34 \pm 0.62 \; (GC)$	3.64 ± 0.62	10.61 ± 0.22
1.0%	$12.69 \pm 0.28 \; (GC)$	5.48 ± 0.42	15.85 ± 0.16
5.0%	25.7 ± 0.55 (GC)	14.70 ± 0.24	31.64 ± 0.19
10.0%	$36.91 \pm 0.18 (GC)$	22.74 ± 1.67	41.32 ± 0.48
20.0%	48.14 ± 0.23 (GC)	51.83 ± 0.52	54.18 ± 0.51
30.0%	56.1 ± 0.32 (GC)	57.79 ± 1.61	60.82 ± 0.08
40.0%	$62.52 \pm 0.23 \; (GC)$	64.92 ± 0.93	65.01 ± 0.24
50.0%	68.62 ± 1.20 (GC)	69.11 ± 0.62	68.67 ± 0.17
60.0%	73.34 ± 0.39 (F)	71.87 ± 0.91	71.10 ± 0.38
Rank	2.33	2.33	1.33

Table 9: Comparison of results for CIFAR-100 on ViT Small

Size	DeepCore	Moderate-DS	Our Method
0.5%	$7.45 \pm 0.17 \; (GC)$	5.20 ± 0.06	11.03 ± 0.23
1.0%	$14.57 \pm 0.47 \; (GC)$	7.82 ± 0.20	14.35 ± 0.18
5.0%	24.50 ± 0.71 (GC)	21.12 ± 0.30	24.85 ± 0.30
10.0%	32.00 ± 0.34 (GC)	27.00 ± 0.50	32.60 ± 0.30
20.0%	40.73 ± 1.19 (GC)	38.51 ± 0.47	41.30 ± 0.80
30.0%	$45.50 \pm 1.01 \; (GC)$	44.90 ± 0.31	48.12 ± 0.48
40.0%	46.72 ± 1.02 (GC)	46.20 ± 0.60	51.75 ± 2.23
50.0%	$47.11 \pm 2.00 \; (GC)$	52.83 ± 0.32	57.62 ± 0.64
60.0%	$49.66 \pm 0.64 \text{ (GC)}$	55.00 ± 0.40	60.81 ± 0.53
Rank	2.11	2.77	1.11

Table 11: Comparison of results for Tiny ImageNet on ResNet-50

Size	DeepCore	Moderate-DS	Our Method
$\begin{array}{c} 0.5\% \\ 1.0\% \\ 5.0\% \\ 10.0\% \\ 20.0\% \\ 30.0\% \\ 40.0\% \\ 60.0\% \end{array}$	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 1.71 \pm 0.27 \\ 3.38 \pm 0.13 \\ 12.39 \pm 0.02 \\ 23.21 \pm 0.23 \\ 34.71 \pm 0.36 \\ 41.65 \pm 0.20 \\ 48.72 \pm 0.61 \\ 56.10 \pm 0.13 \end{array}$	$\begin{array}{c} 5.24\pm 0.03\\ 10.55\pm 0.41\\ 26.07\pm 0.19\\ 34.08\pm 0.21\\ 43.95\pm 1.03\\ 50.82\pm 0.03\\ 55.23\pm 0.08\\ 61.38\pm 0.11\end{array}$
Rank	2.00	3.00	1.00

Table 13: Comparison of results for ImageNet-1K on ResNet-18

Size	DeepCore	Noise-free Gradients
$\begin{array}{c} 0.1\% \\ 0.5\% \\ 1.0\% \\ 5.0\% \\ 10.0\% \end{array}$	$\begin{array}{c} 1.29 \pm 0.09 \; ({\rm CAL}) \\ 7.66 \pm 0.43 \; ({\rm GraphCut}) \\ 18.10 \pm 0.22 \; ({\rm GraNd}) \\ \textbf{47.64 \pm 0.03 \; (Forgetting)} \\ \textbf{55.12 \pm 0.13 \; (Forgetting)} \end{array}$	$\begin{array}{c} {\bf 1.88 \pm 0.10} \\ {\bf 11.58 \pm 0.15} \\ {\bf 22.82 \pm 0.07} \\ {\bf 43.50 \pm 0.41} \\ {\bf 49.51 \pm 0.26} \end{array}$
Rank	1.6	1.4

We visualize the top and bottom-ranked samples in this subsection according to our coreset selection algorithm. We chose the "Zebra" class from the 1000 ImageNet classes. Figure 3 presents the top-ranked and bottom-ranked 12 images of the "Zebra" class. The selected samples clearly emphasize the fact that the images with top ranks are representative of the class. In contrast, images at the bottom positions are ambiguously labelled (have parts of objects from other classes) or difficult to label as "Zebra". As our objective is to select the most representative images from each class, the visualization provides evidence for the intuition over which the proposed method is founded. We have provided additional examples in the supplementary material.



Top ranked images

Bottom ranked images

Figure 3: Top-ranked and bottom-ranked 12 images from the "zebra" class from ImageNet-1K dataset by the proposed "Noise-free Gradients" approach. Images are ranked as per the number of neighbors within a threshold gradient similarity value. Top-ranked images are unambiguous representative of the class, while bottom-ranked images are either mislabeled or ambiguously labeled.

5.2.4 Robustness against image noise

In realistic scenarios, training data may be polluted by corrupted images (Xia et al., 2023). We have randomly added five different types of noise to the training subset of the original dataset, namely Gaussian Noise, random occlusion, resolution change, fog and motion blur. We have conducted coreset selection on the corrupted train set and observed its generalization performance on the uncorrupted test set.

We have adapted the implementation provided by Moderate-DS(Xia et al., 2023) available in public github repository.⁵ Table 14 compares accuracy obtained by all the methods available in DeepCore library and Moderate with our *Noise-free Gradients* method for a 30% of the training dataset impacted by noise. ResNet-18 architecture is used for evaluation. As can be seen, for this significant amount of noise added to the original dataset, our method outperforms all the other methods for most of the coreset selection. As our proposed method composes the coreset with images having gradient similarity with a higher number of images from the same class, the coreset is able to perform better than existing methods even when the training dataset is corrupted by significant image noise.

Table 14:	Comparison of results for CIFAR-100 with 30 ^o	%
	corruption on ResNet-18	

Table 15:	Comparison	of	results	for	CIFAR-100	with	5%
	corrup	tio	on on R	esN	et-50		

Size	DeepCore	Moderate-DS	Our Method
$\begin{array}{c} 0.5\% \\ 1.0\% \\ 5.0\% \\ 10.0\% \\ 20.0\% \\ 30.0\% \end{array}$	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 5.07 \pm 0.57 \\ 17.56 \pm 0.61 \\ 19.03 \pm 0.8 \\ 29.20 \pm 0.23 \\ 41.27 \pm 0.56 \\ \textbf{57.58} \pm \textbf{0.16} \end{array}$	$\begin{array}{c} 11.1 \pm 0.18 \\ 33.56 \pm 0.65 \\ 23.65 \pm 0.15 \\ 31.24 \pm 0.11 \\ 41.89 \pm 0.40 \\ 51.93 \pm 0.25 \end{array}$
Rank	2.83	2	1.17

Size	DeepCore	Moderate-DS	Our Method
$10.0\% \\ 20.0\% \\ 30.0\% \\ 40.0\% \\ 60.0\%$	$ \begin{vmatrix} 22.2 \pm 0.13 & (H) \\ 42.50 \pm 1.27 & (H) \\ 53.88 \pm 0.37 & (H) \\ 60.54 \pm 0.94 & (H) \\ 70.22 \pm 0.22 & (F) \end{vmatrix} $	$\begin{array}{c} 35.1 \pm 0.34 \\ 46.78 \pm 1.9 \\ 57.36 \pm 1.22 \\ \textbf{65.4} \pm \textbf{0.39} \\ \textbf{71.46} \pm \textbf{0.19} \end{array}$	$\begin{array}{c} \textbf{42.0} \pm \textbf{0.24} \\ \textbf{53.48} \pm \textbf{0.22} \\ \textbf{59.9} \pm \textbf{0.13} \\ \textbf{64.96} \pm \textbf{0.2} \\ \textbf{70.27} \pm \textbf{0.16} \end{array}$
Rank	3.00	1.60	1.40

Table 15 and Table 16 show accuracy comparison for 5% and 20% of images from CIFAR-100 dataset impacted by noises, trained on ResNet-50 architecture. Table 17 compares accuracies obtained for 20% noise

⁵https://github.com/tmllab/2023_ICLR_Moderate-DS

Size	DeepCore	Moderate-DS	Our Method
10.0%	24.78 ± 0.45 (F)	27.05 ± 0.8	42.32 ± 0.23
20.0%	44.42 ± 0.46 (H)	42.98 ± 0.87	54.63 ± 0.26
30.0%	53.57 ± 0.31 (H)	55.80 ± 0.96	60.58 ± 0.20
40.0%	60.72 ± 1.78 (H)	61.84 ± 1.96	64.83 ± 0.35
60.0%	69.10 ± 1.73 (H)	70.05 ± 1.29	70.51 ± 1.31
Rank	2.80	2.20	1.00

Table 16: Comparison of results for CIFAR-100 with 20% corruption on ResNet-50

ratio in the CIFAR-100 dataset, trained on transformer architecture(ViT Small is used). In all four studies, we can see that our method significantly outperforms all the SOTA methods.

5.2.5 Robustness against label noise

Practical datasets may also involve label noise, where some images are mislabeled. We have adapted the implementation provided by Moderate-DS (Xia et al., 2023) available in public GitHub repository⁶ to generate mislabeled data. Table 18 and Table 19 compare accuracy obtained by all the methods available in the DeepCore library and Moderate with our *Noise-free Gradients* method for a 20% of the training dataset impacted by label noise on ResNet-50 architecture and ViT Small architecture respectively. As can be seen, for this significant amount of noise added to the original dataset labels, our method outperforms all the other methods. As our proposed method composes the coreset with images having gradient similarity with a higher number of images from the same class, the coreset is able to perform better than existing methods even when the training dataset is corrupted by significant label noise.

Table 18: Comparison of results on CIFAR-100 Dataset with 20% label noise on ResNet-50 architecture

Table 19:	Comparison of	results on	CIFAR-100	Dataset	with 20%
	label noise	on ViT Si	mall architec	ture	

Table 17: Comparison of results for CIFAR-100 with 20%

corruption on ViT Small

DeepCore	Moderate-DS	Our Method	Size	DeepCore	Moderate-DS	Our Method
$38.59 \pm 1.29 (GC)$ $47.15 \pm 1.66 (GC)$ $57.04 \pm 0.88 (GC)$	25.02 ± 0.27 42.75 ± 1.20 56.28 ± 1.11	$\begin{array}{c} 40.71 \pm 0.34 \\ 52.99 \pm 0.51 \\ 60.77 \pm 1.17 \end{array}$	10.0% 20.0% 30.0%	$ \begin{vmatrix} 28.66 \pm 0.52 & (GC) \\ 38.79 \pm 0.55 & (GC) \\ 40.22 \pm 0.92 & (GC) \end{vmatrix} $	$\begin{array}{c} 23.02 \ \pm 0.32 \\ 34.92 \ \pm 0.10 \\ 43.91 \ \pm 0.28 \end{array}$	$\begin{array}{c} 31.39\pm0.36\ 36.90\pm0.21\ 44.71\pm0.07 \end{array}$
2.00	3.00	1.00	Rank	2.00	3.00	1.00

5.2.6 Impact on class-wise accuracy

Size | 10.0% | 20.0% | 30.0% | Rank |

We visualize the class-wise accuracy of training a ResNet-18 model on the CIFAR-100 dataset, comparing the performance of the full dataset and coreset selection with 20% of the dataset in Figure 4. We observe a Spearman rank-order correlation coefficient of 0.87, indicating a high correlation between class-wise accuracies obtained with the full dataset and 20% of the full dataset as coreset. It can be seen that our method does not adversely impact any particular class during coreset selection. The coreset composed by our method results in class-wise accuracies strongly correlated to that of the full dataset.

5.2.7 Cross-architecture generalization

To study the robustness of our coreset selection method, we have carried out cross-architecture generalization performance comparison at 1% and 10% fractions of CIFAR-10 and Tiny ImageNet for two CNN architectures ResNet-18 (He et al., 2016), VGG-16 (Simonyan & Zisserman, 2015) and one transformer based architecture ViT-Small. The comparative analysis is tabulated in Table 20 and Table 21. The "Source" column denotes the architecture on which the coreset is selected, and the "Target" row indicates the architecture on which its performance is evaluated. As can be observed, our proposed method is able to perform consistently better than existing methods in a cross-architecture setting. Cross-architecture analysis on the Tiny-ImageNet dataset is provided in the supplementary material.

⁶https://github.com/tmllab/2023_ICLR_Moderate-DS



Figure 4: Class-wise accuracy comparison between full dataset and coreset with 20% selection percentage. A high correlation is observed indicating coreset composed by our proposed method not impacting any particular class adversely.

Table 20:	Cross-architecture comparison	for	1%	coreset o	of
	CIFAR-10				

Table 21: Cross-architecture comparison for 10% coreset of the CIFAR-10 $\,$

Target \rightarrow		ResNet-18	
Source \downarrow	DeepCore	Moderate-DS	Our Method
ResNet-18 VGG-16 ViT Small	$\begin{array}{c} 42.78 \pm 1.30 \; ({\rm GC}) \\ 43.02 \pm 1.30 \; ({\rm GC}) \\ 26.01 \pm 2.00 \; ({\rm GC}) \end{array}$	$\begin{array}{r} 41.44 \pm 0.34 \\ 42.12 \pm 0.27 \\ 41.76 \pm 0.35 \end{array}$	$48.00 \pm 2.10 \\ 46.27 \pm 0.33 \\ 45.05 \pm 0.51$
Target \rightarrow		VGG-16	
Source \downarrow	DeepCore	Moderate-DS	Our Method
ResNet-18 VGG-16 ViT Small	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 44.84 \pm 0.33 \\ 44.35 \pm 0.45 \\ 34.44 \pm 0.30 \end{array}$	$47.21 \pm 0.95 \\ 47.64 \pm 0.71 \\ 38.43 \pm 0.40$
Target \rightarrow		ViT Small	
Source \downarrow	DeepCore	Moderate-DS	Our Method
ResNet-18 VGG-16 ViT Small	$\begin{array}{c} 29.06 \pm 0.75 \; (\text{GC}) \\ 42.06 \pm 0.90 \; (\text{GC}) \\ 22.89 \pm 1.45 \; (\text{GC}) \end{array}$	$\begin{array}{r} 34.71 \pm 0.23 \\ 31.25 \pm 0.73 \\ 34.44 \pm 0.30 \end{array}$	$40.35 \pm 0.90 \\ 44.89 \pm 0.55 \\ 38.43 \pm 0.40$

Target \rightarrow		ResNet-18	
Source \downarrow	DeepCore	Moderate-DS	Our Method
ResNet-18 VGG-16 ViT Small	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 75.68 \pm 0.62 \\ 74.27 \pm 0.72 \\ 73.66 \pm 0.29 \end{array}$	$\begin{array}{c} 78.62 \pm 1.05 \\ 78.75 \pm 0.16 \\ 75.15 \pm 0.42 \end{array}$
Target \rightarrow		VGG-16	
Source \downarrow	DeepCore	Moderate-DS	Our Method
ResNet-18 VGG-16 ViT Small	$ \begin{vmatrix} 75.29 \pm 1.05 & (GC) \\ 77.91 \pm 0.71 & (GC) \\ 70.38 \pm 0.87 & (GC) \end{vmatrix} $	75.60 ± 0.86 76.62 ± 0.52 73.81 ± 0.68	$78.22 \pm 0.29 \\78.84 \pm 0.33 \\75.47 \pm 0.35$
Target \rightarrow		ViT Small	
Source \downarrow	DeepCore	Moderate-DS	Our Method
ResNet-18 VGG-16 ViT Small	$ \begin{array}{c} 62.99 \pm 0.12 \; (\text{GC}) \\ 63.14 \pm 0.18 \; (\text{GC}) \\ 49.30 \pm 1.30 \; (\text{GC}) \end{array} $	$\begin{array}{c} 61.69 \pm 0.23 \\ 61.49 \pm 0.33 \\ 54.40 \pm 2.35 \end{array}$	$63.66 \pm 0.35 \\ 66.88 \pm 0.41 \\ 58.32 \pm 0.16$

5.3 GAN Analysis

So far, we have presented the effectiveness of our proposed method in the context of training a discriminative model (classifier); however, in this subsection, we investigate the effectiveness of the proposed method in a different context, which is training a generative model. Similar to our approach for the classification task, we train the generative model for the initial few epochs with the full dataset and then select the class-wise coreset with a 50% selection percentage. We utilize the discriminator as a classifier for the purpose of computation of gradient similarity and ranking of images. Thereon, the generative model is trained with the coreset.

5.3.1 Results on CIFAR-10

We have trained four different GAN architectures, namely DCGAN, MSGAN, SAGAN, and SNGAN, on the CIFAR-10 dataset. Table 22 compared the FID and IS scores obtained by training with a full dataset and coreset composed by small GAN, instance selection, and our method. As can be observed, our method consistently performs better than the other two methods in the FID score, which is a key metric for measuring the quality of generated images.

Method	DCGAN		MSGAN		SAGAN		SNGAN	
	FID \downarrow	IS \uparrow	$\mathrm{FID}\downarrow$	IS \uparrow	$\mathrm{FID}\downarrow$	IS \uparrow	$\mathrm{FID}\downarrow$	IS \uparrow
Full Dataset	28.80	5.84	33.16	5.81	33.73	5.56	22.44	7.36
Small GAN	35.53	5.59	45.76	5.28	44.97	4.92	35.89	6.34
Instance Selection	31.26	5.62	39.05	5.71	36.96	5.48	27.81	6.98
Our Method	31.14	5.70	35.65	5.51	34.60	5.42	24.46	7.22

Table 22: Comparison of results on CIFAR-10 Dataset. \downarrow means lower is better, \uparrow means higher is better

5.3.2 Results on CIFAR-100

We have trained four different GAN architectures, namely DCGAN, MSGAN, SAGAN and SNGAN, on the CIFAR-100 dataset. Table 23 compared the FID and IS scores obtained by training with full dataset and coreset composed by small GAN, instance selection and our method. As can be observed, our method consistently performs better than the other two methods in FID score, which is a key metric for measuring the quality of generated images.

Table 23: Comparison of results on CIFAR-100 Dataset. ↓ means lower is better, ↑ means higher is better

Method	DCGAN		MSGAN		SAGAN		SNGAN	
	FID \downarrow	IS \uparrow	$\mathrm{FID}\downarrow$	IS \uparrow	$\mathrm{FID}\downarrow$	IS \uparrow	$\mathrm{FID}\downarrow$	IS \uparrow
Full Dataset	27.98	6.13	37.59	7.26	33.37	5.79	26.97	6.26
Small GAN	32.81	6.11	71.13	6.75	41.53	5.76	39.71	5.51
Instance Selection	30.48	6.07	53.88	7.10	39.17	5.60	28.88	6.13
Our Method	30.08	6.08	50.22	6.64	36.15	5.61	27.71	6.14

5.3.3 Results on Tiny ImageNet

We have trained SAGAN architecture on the Tiny ImageNet dataset. Table 24 compared the FID and IS scores obtained by training with full dataset and coreset composed by small GAN, instance selection and our method.

Table 24: Comparison of results on Tiny ImageNet Dataset. ↓ means lower is better, ↑ means higher is better

Method	SAGAN		
	$\mathrm{FID}\downarrow$	IS \uparrow	
Full Dataset	55.67	7.12	
Small GAN	75.12	6.65	
Instance Selection	69.16	6.46	
Our Method	56.81	6.97	

5.4 Computation time analysis

5.4.1 Image Classification

Our method does not consider all images in a given dataset in a single coreset selection step. Instead, it processes the dataset class-wise, implements multiprocessing and optimizes the nearest neighbour algorithm to parallelize the ranking process. Hence, the methodology can be easily scaled to a large-scale dataset beyond ImageNet-1K. Also, the computation of gradients and nearest neighbours is a one-time process for a given dataset and architecture (per selection epoch). After that, aggregating the individual image ranks across the training epochs takes only 225 seconds (steps 11-13 in Algorithm 1). As we directly utilize the gradients of the loss function with respect to the weights of the last fully connected layer, there is no requirement for storage of the gradient values. Table 25 provides timing analysis for the ImageNet dataset on ResNet-18 architecture on a 48 GB RTX A6000 GPU.

Process	Execution Time
Single epoch full dataset training	4.46 hours
Gradient Calculation (per epoch)	3.19 hours
Ranking for the coreset	225 seconds

5.4.2 GAN training

By using 50% of the total dataset, we train the GAN model with coreset selection for half the number of update steps compared to training with the entire dataset. As shown in section 5.3, our method can achieve a comparable FID score to training with a whole dataset while utilizing half of the available dataset. Table 26 provides the total training time required for various methods, and we can see that our method can reduce the total time taken significantly. At the same time, degradation in performance is not significant.

Table 26: Execution time comparison for GAN training on CIFAR-100 Dataset (in hours)

Method	DCGAN	MSGAN	SAGAN	SNGAN
Full Dataset Small GAN Instance Selection Our Method	2.47 9.82 1.32 1.34	$5.14 \\ 15.10 \\ 3.00 \\ 3.53$	$2.90 \\ 16.22 \\ 1.56 \\ 1.80$	9.48 19.35 3.95 3.98

5.5 Accuracy vs training speed tradeoff

Figure 5 and 6 shows classification accuracy vs computational time trade-off for full dataset training vs coreset training using 10%, 20% and 30% of the entire dataset. As it can be seen, with a small reduction in accuracy, we can improve the training speed, and it shows the effectiveness of the proposed coreset method during hyper-parameter tuning, which may require evaluating different values for them over multiple runs.



Figure 5: Accuracy vs training speed trade-off. It can be observed that using 30% of coreset results in similar accuracy while the training time is reduced more than 3 times.



Figure 6: Accuracy vs training speed trade-off. For a complex dataset like Tiny ImageNet, the reduction in training time is not insignificant for applications like hyperparameter tuning where model needs to be trained multiple times.

6 Conclusions

In this paper, we introduced *Noise-free Gradients*, an intuition-driven, gradient similarity-based coreset selection method for identifying representative instances from large training datasets. We studied its effectiveness through our extensive experiments on popular object recognition datasets (CIFAR-10, CIFAR-100, Tiny ImageNet, ImageNet-1K) and sophisticated model architectures (CNNs and Vision Transformer) at various coreset sizes. We demonstrated superior generalization performance of classifiers trained on the resulting coresets. We have also demonstrated that our model achieves consistently higher performance in the case of cross-architecture generalization. Thus, we strongly propose our method as an essential baseline to benchmark the forthcoming coreset selection methods.

A branch of existing methods, such as Forgetting (Toneva et al., 2019), utilizes noisy gradients to compose coresets. The argument put forward by these methods is the samples that are hard to learn will be able to provide better generalization on the test set. However, our approach explores an opposing notion. Samples that derive loss gradients similar to many other samples can effectively represent the whole dataset and thus form a coreset. This hypothesis naturally exploits the redundancy in large datasets to compose extremely small coresets. Our method outperforms other coreset selection methods even in the presence of noise in the original dataset. We have also applied our coreset method to improve computational time significantly for GAN training without sacrificing much performance.

6.1 Future Directions

While one can appreciate the simplicity and intuition behind the proposed coreset selection method, it must be studied further, particularly concerning the approximation error between the gradients computed by the coreset and the entire dataset. Inducing diversity in the selected samples along with individual importance weights (or per-sample step size) would further improve the effectiveness of the proposed approach. One may consider these aspects for future study.

There have been a few reported applications of coreset selection, namely Active Learning (Sener & Savarese, 2018) and semi-supervised learning (Killamsetty et al., 2021c). It would be interesting to study the effectiveness of our *Noise-free Gradients* approach concerning these application areas.

References

- Sharat Agarwal, Himanshu Arora, Saket Anand, and Chetan Arora. Contextual diversity for active learning, 2020.
- Jordan T. Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep batch active learning by diverse, uncertain gradient lower bounds, 2020.
- Cody Coleman, Christopher Yeh, Stephen Mussmann, Baharan Mirzasoleiman, Peter Bailis, Percy Liang, Jure Leskovec, and Matei Zaharia. Selection via proxy: Efficient data selection for deep learning, 2020.
- Terrance DeVries, Michal Drozdzal, and Graham W. Taylor. Instance selection for gans, 2020. URL https://arxiv.org/abs/2007.15255.
- Melanie Ducoffe and Frederic Precioso. Adversarial active learning for deep networks: a margin based approach, 2018.
- Dan Feldman. Introduction to core-sets: an updated survey, 2020.
- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014. URL https://arxiv.org/abs/ 1406.2661.
- Chengcheng Guo, Bo Zhao, and Yanbing Bai. Deepcore: A comprehensive library for coreset selection in deep learning, 2022.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In Proceedings of 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR '16, pp. 770-778. IEEE, June 2016. doi: 10.1109/CVPR.2016.90. URL http://ieeexplore.ieee.org/document/ 7780459.
- Melissa Heikkilä. We're getting a better idea of AI's true carbon footprint technologyreview.com. https://www.technologyreview.com/2022/11/14/1063192/ were-getting-a-better-idea-of-ais-true-carbon-footprint?truid=&utm_source=the_ algorithm&utm_medium=email&utm_campaign=the_algorithm.unpaid.engagement&utm_content= 12-04-2023.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium, 2018. URL https://arxiv.org/abs/1706.08500.
- Rishabh Iyer, Ninad Khargonkar, Jeff Bilmes, and Himanshu Asnani. Submodular combinatorial information measures with applications in machine learning, 2021.
- Angelos Katharopoulos and François Fleuret. Not all samples are created equal: Deep learning with importance sampling, 2019.
- Krishnateja Killamsetty, Durga Sivasubramanian, Ganesh Ramakrishnan, Abir De, and Rishabh Iyer. Gradmatch: Gradient matching based data subset selection for efficient deep model training, 2021a.
- Krishnateja Killamsetty, Durga Sivasubramanian, Ganesh Ramakrishnan, and Rishabh Iyer. Glister: Generalization based data subset selection for efficient and robust learning, 2021b.
- Krishnateja Killamsetty, Xujiang Zhao, Feng Chen, and Rishabh Iyer. Retrieve: Coreset selection for efficient and robust semi-supervised learning. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), Advances in Neural Information Processing Systems, volume 34, pp. 14488-14501. Curran Associates, Inc., 2021c. URL https://proceedings.neurips.cc/paper_files/ paper/2021/file/793bc52a941b3951dfdb85fb04f9fd06-Paper.pdf.

- Krishnateja Killamsetty, Xujiang Zhao, Feng Chen, and Rishabh Iyer. Retrieve: Coreset selection for efficient and robust semi-supervised learning, 2021d.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. pp. 32-33, 2009. URL https: //www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf.
- Ya Le and Xuan S. Yang. Tiny imagenet visual recognition challenge. 2015. URL https://api. semanticscholar.org/CorpusID:16664790.
- Seung Hoon Lee, Seunghyun Lee, and Byung Cheol Song. Vision transformer for small-size datasets, 2021.
- Siwen Liu. A Review for Submodular Optimization on Machine Scheduling Problems, pp. 252–267. Springer International Publishing, Cham, 2020. ISBN 978-3-030-41672-0. doi: 10.1007/978-3-030-41672-0_16. URL https://doi.org/10.1007/978-3-030-41672-0_16.
- Qi Mao, Hsin-Ying Lee, Hung-Yu Tseng, Siwei Ma, and Ming-Hsuan Yang. Mode seeking generative adversarial networks for diverse image synthesis, 2019. URL https://arxiv.org/abs/1903.05628.
- Katerina Margatina, Giorgos Vernikos, Loïc Barrault, and Nikolaos Aletras. Active learning by acquiring contrastive examples, 2021.
- Baharan Mirzasoleiman, Jeff Bilmes, and Jure Leskovec. Coresets for data-efficient training of machine learning models, 2020.
- Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks, 2018. URL https://arxiv.org/abs/1802.05957.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, highperformance deep learning library. In Advances in Neural Information Processing Systems 32, pp. 8024-8035. Curran Associates, Inc., 2019. URL http://papers.neurips.cc/paper/ 9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf.
- Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet: Finding important examples early in training, 2023.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. Journal of machine learning research, 12(Oct):2825–2830, 2011.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge, 2015.
- Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans, 2016. URL https://arxiv.org/abs/1606.03498.

Roy Schwartz, Jesse Dodge, Noah A. Smith, and Oren Etzioni. Green ai, 2019.

- Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach, 2018.
- Burr Settles. Active Learning. Springer International Publishing, 2012. ISBN 9783031015601. doi: 10.1007/978-3-031-01560-1. URL http://dx.doi.org/10.1007/978-3-031-01560-1.
- Yanyao Shen, Hyokun Yun, Zachary C. Lipton, Yakov Kronrod, and Animashree Anandkumar. Deep active learning for named entity recognition, 2018.

- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2015.
- Samarth Sinha, Han Zhang, Anirudh Goyal, Yoshua Bengio, Hugo Larochelle, and Augustus Odena. Smallgan: Speeding up gan training using core-sets, 2019.
- Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J. Gordon. An empirical study of example forgetting during deep neural network learning, 2019.
- Max Welling. Herding dynamic weights for partially observed random field models, 2012.
- Xiaobo Xia, Jiale Liu, Jun Yu, Xu Shen, Bo Han, and Tongliang Liu. Moderate coreset: A universal method of data selection for real-world data-efficient deep learning. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=7D5EECbOaf9.
- Shuo Yang, Zeke Xie, Hanyu Peng, Min Xu, Mingming Sun, and Ping Li. Dataset pruning: Reducing training data by examining generalization influence, 2023.
- Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks, 2019. URL https://arxiv.org/abs/1805.08318.