# Reinforcement Learning-Guided Data Selection via Redundancy Assessment

## Supplementary Material

## 1. Proof of Proposition 1

**Proof 1** *For simplicity, we approximate the fully connected layer as a linear model, i.e., $\hat{\boldsymbol{y}}_i = \boldsymbol{w}\boldsymbol{x}_i + \boldsymbol{b}$.*

$$\left\| \hat{\boldsymbol{y}}_i - \hat{\boldsymbol{y}}_j \right\| = \left\| \boldsymbol{w}(\boldsymbol{x}_i - \boldsymbol{x}_j) \right\| \tag{12}$$
$$\leq \left\| \boldsymbol{w} \right\| \left\| \boldsymbol{x}_i - \boldsymbol{x}_j \right\| \tag{13}$$
$$\leq \epsilon \left\| \boldsymbol{w} \right\| \tag{14}$$
$$\tag{15}$$

*In the above, the Inequality (13) follows from Hölder's inequality. For a given model with corresponding $\boldsymbol{w}$, $\left\| \hat{\boldsymbol{y}}_i - \hat{\boldsymbol{y}}_j \right\| \leq \mathcal{O}(\epsilon)$ holds.*

## 2. Proof of Proposition 2

**Proof 2** *Let $\boldsymbol{x}_i$ and $\boldsymbol{x}_j$ are two samples, and $\boldsymbol{x}_i$ is $\epsilon$-covered by $\boldsymbol{x}_j$. The cross-entropy loss of $\boldsymbol{x}_i$ is given as*

$$loss(\boldsymbol{x}_i) = L(\boldsymbol{w}\boldsymbol{x}_i + b) = -\sum_{c=1}^{K} y_{ic} \log(\boldsymbol{w}\boldsymbol{x}_i + b)_c. \tag{16}$$

*Let $\hat{y}_{ic} = (\boldsymbol{w}\boldsymbol{x}_i + b)_c$, then $loss(\boldsymbol{x}_i) - loss(x_j)$ can be computed as*

$$loss(\boldsymbol{x}_i) - loss(x_j) = L(\boldsymbol{w}\boldsymbol{x}_i + b) - L(\boldsymbol{w}\boldsymbol{x}_j + b)$$
$$= -\sum_{c=1}^{K} y_{ic}(\log(\hat{y}_{ic}) - \log(\hat{y}_{jc}))$$
$$= -\sum_{c=1}^{K} y_{ic} \log \frac{\hat{y}_{ic}}{\hat{y}_{jc}}$$
$$= -log \frac{\hat{y}_{ik}}{\hat{y}_{jk}}. \tag{17}$$

*According to Definition 1, $\hat{y}_{ik}/\hat{y}_{jk}$ can be re-written as follows,*

$$\frac{\hat{y}_{ik}}{\hat{y}_{jk}} = \frac{(\boldsymbol{w}\boldsymbol{x}_i + b)_k}{(\boldsymbol{w}\boldsymbol{x}_j + b)_k} \leq \frac{(\boldsymbol{w}(\boldsymbol{x}_i + \epsilon) + b)_k}{(\boldsymbol{w}\boldsymbol{x}_j + b)_k}. \tag{18}$$

*Substituting Eq. (18) into Eq. (17), the gap in losses can be derived as follows.*

$$\lim_{\epsilon \to 0} \frac{\hat{y}_{ik}}{\hat{y}_{jk}} \leq \lim_{\epsilon \to 0} \frac{(\boldsymbol{w}(\boldsymbol{x}_i + \epsilon) + b)_k}{(\boldsymbol{w}\boldsymbol{x}_j + b)_k} = \lim_{\epsilon \to 0} \frac{(\boldsymbol{w}\boldsymbol{x}_j + b)_k + \epsilon w_k}{(\boldsymbol{w}\boldsymbol{x}_j + b)_k} = 1. \tag{19}$$

*Similarly,*

$$\lim_{\epsilon \to 0} \frac{\hat{y}_{ik}}{\hat{y}_{jk}} \geq \lim_{\epsilon \to 0} \frac{(\boldsymbol{w}(\boldsymbol{x}_i - \epsilon) + b)_k}{(\boldsymbol{w}\boldsymbol{x}_j + b)_k} = 1. \tag{20}$$

*Thus, according to Eq. (19) and Eq. (20), $\lim_{\epsilon \to 0} \frac{\hat{y}_{ik}}{\hat{y}_{jk}} = 1$, and we have*

$$\lim_{\epsilon \to 0} loss(\boldsymbol{x}_i) - loss(\boldsymbol{x}_j) = 0. \tag{21}$$

## 3. Proof of Lemma 1

**Proof 3** *Without loss of generality, the gradient of loss for $\boldsymbol{x}$ w.r.t. the weight $W_h$ of the $h$-th intermediate layer is given as*

$$g_{W_h} = (\frac{\partial}{\partial W_h^T} L(f_\theta(\boldsymbol{x}), y))^T. \tag{22}$$

*According to the chain rule, the gradient of the $i$-th element of the weight $W_h$ can be written as*

$$g_{W_{hi}^T} = \frac{\partial L(f_\theta(\boldsymbol{x}), y)}{\partial \tilde{\boldsymbol{x}}_{h+1}^T} \frac{\partial \tilde{\boldsymbol{x}}_{h+1}^T}{\partial W_{hi}^T}$$
$$= \frac{\partial L(f_\theta(\boldsymbol{x}), y)}{\partial \tilde{\boldsymbol{x}}_{h+1}^T} \tilde{\boldsymbol{x}}_h^T, \tag{23}$$

*where $\tilde{\boldsymbol{x}}_{h+1}$ is the output obtained from the $h$-th intermediate layer, i.e., $\tilde{\boldsymbol{x}}_{h+1} = W_h^T \tilde{\boldsymbol{x}}_h + b$, where $\tilde{\boldsymbol{x}}_h$ can be understood as the feature map from the intermediate layer. Then, we have*

$$g_{W_h} = [g_{W_{h1}^T}, g_{W_{h2}^T}, ..., g_{W_{hD}^T}]^T = [\frac{\partial L(f_\theta(\boldsymbol{x}), y)}{\partial \tilde{\boldsymbol{x}}_{h+1}^T} \tilde{\boldsymbol{x}}_h^T]^T$$
$$= \tilde{\boldsymbol{x}}_h \frac{\partial L(f_\theta(\boldsymbol{x}), y)}{\partial \tilde{\boldsymbol{x}}_{h+1}^T} = \tilde{\boldsymbol{x}}_h g_{\tilde{\boldsymbol{x}}}^T \tag{24}$$

*According to Definition 1, we have $\left\| \tilde{\boldsymbol{x}}_{i,h} - \tilde{\boldsymbol{x}}_{j,h} \right\| \leq \epsilon$. Let $\Delta \boldsymbol{x} = \tilde{\boldsymbol{x}}_{i,h} - \tilde{\boldsymbol{x}}_{j,h}$, i.e., $\Delta \boldsymbol{x} \in \mathbb{R}^D$ and $\left\| \Delta \boldsymbol{x} \right\| \leq \epsilon$. Meanwhile, we can derive the following equation,*

$$\tilde{\boldsymbol{x}}_{i,h+1} = W_h^T \tilde{\boldsymbol{x}}_{i,h} + b = \tilde{\boldsymbol{x}}_{j,h+1} + W_h^T \Delta \boldsymbol{x} \tag{25}$$

*According to Eq. (24) and Eq. (25), for $\boldsymbol{x}_i$, the gradient of loss w.r.t. the weight $W_h$ is given as*

$$g_{W_h}^{\boldsymbol{x}_i} = (\tilde{\boldsymbol{x}}_{j,h} + \Delta \boldsymbol{x}) g_{\tilde{\boldsymbol{x}} + \Delta \boldsymbol{x}}^T. \tag{26}$$

*We further use the first-order Taylor expansion to decompose the gradient $g_{\tilde{\boldsymbol{x}} + \Delta \boldsymbol{x}}$*

$$g_{\tilde{\boldsymbol{x}} + \Delta \boldsymbol{x}} = g_{\tilde{\boldsymbol{x}}} + H_{\tilde{\boldsymbol{x}}} \Delta \boldsymbol{x} + \mathcal{R}_2(\Delta \boldsymbol{x}), \tag{27}$$

*where $\mathcal{R}_2(\Delta \boldsymbol{x})$ denotes the terms no less than the second order. Substituting Eq. (27) back to Eq. (26), the gradient for parameter update is*

$$g_{W_h}^{\boldsymbol{x}_i} = (\tilde{\boldsymbol{x}}_{j,h} + \Delta \boldsymbol{x})(g_{\tilde{\boldsymbol{x}}} + H_{\tilde{\boldsymbol{x}}} \Delta \boldsymbol{x} + \mathcal{R}_2(\Delta \boldsymbol{x}))^T \tag{28}$$

*In this way,*

$$
\begin{aligned}
\Delta g_{W_h} &= g_{W_h}^{\boldsymbol{x}_i} - g_{W_h}^{\boldsymbol{x}_j} \\
&= \tilde{\boldsymbol{x}}_{j,h}(H_{\tilde{\boldsymbol{x}}}\Delta\boldsymbol{x})^T + \Delta\boldsymbol{x}(g_{\tilde{\boldsymbol{x}}} + H_{\tilde{\boldsymbol{x}}}\Delta\boldsymbol{x})^T \\
&\quad + (\tilde{\boldsymbol{x}}_{j,h} + \Delta\boldsymbol{x})\mathcal{R}_2(\Delta\boldsymbol{x})^T \\
&\approx \tilde{\boldsymbol{x}}_{j,h}(H_{\tilde{\boldsymbol{x}}}\Delta\boldsymbol{x})^T + \Delta\boldsymbol{x}(g_{\tilde{\boldsymbol{x}}} + H_{\tilde{\boldsymbol{x}}}\Delta\boldsymbol{x})^T \\
&= (\tilde{\boldsymbol{x}}_{j,h} + \Delta\boldsymbol{x})(H_{\tilde{\boldsymbol{x}}}\Delta\boldsymbol{x})^T + \Delta\boldsymbol{x}g_{\tilde{\boldsymbol{x}}}^T \\
&= \tilde{\boldsymbol{x}}_{i,h}\Delta\boldsymbol{x}^T H_{\tilde{\boldsymbol{x}}}^T + \Delta\boldsymbol{x}g_{\tilde{\boldsymbol{x}}}^T
\end{aligned}
\tag{29}
$$

## 4. Proof of Proposition 3

**Proof 4** *According to Eq.* (29)*, we have*

$$
\|\Delta g_{W_h}\| = \left\|\tilde{\boldsymbol{x}}_{i,h}\Delta\boldsymbol{x}^T H_{\tilde{\boldsymbol{x}}}^T + \Delta\boldsymbol{x}g_{\tilde{\boldsymbol{x}}}^T\right\| \tag{30}
$$
$$
\leq \left\|\tilde{\boldsymbol{x}}_{i,h}\Delta\boldsymbol{x}^T H_{\tilde{\boldsymbol{x}}}^T\right\| + \left\|\Delta\boldsymbol{x}g_{\tilde{\boldsymbol{x}}}^T\right\| \tag{31}
$$
$$
\leq \|\Delta\boldsymbol{x}\| \left(\|\tilde{\boldsymbol{x}}_{i,h}\|\,\|H_{\tilde{\boldsymbol{x}}}\| + \|g_{\tilde{\boldsymbol{x}}}\|\right). \tag{32}
$$

*Given a model $f_\theta$, since $\|\Delta\boldsymbol{x}\| \leq \epsilon$, the change of the gradient can be formulated as follows,*

$$
\lim_{\epsilon\to 0}\|\Delta g_{W_h}\| = 0. \tag{33}
$$

## 5. Proof of Proposition 4

**Proof 5** *According to Lemma 2, $\forall \boldsymbol{x}' \in -\hat{\mathcal{D}}$, we can find a $\boldsymbol{x} \in \mathcal{D}$ such that $\boldsymbol{x}'$ is $\epsilon$-covered by $\boldsymbol{x}$. The gap in the parameter changes between $\boldsymbol{x}'$ and $\boldsymbol{x}$ can be estimated as follows*

$$
\Delta\mathcal{I}_{up,params} = \mathcal{I}_{up,params}(\boldsymbol{x}) - \mathcal{I}_{up,params}(\boldsymbol{x}') \tag{34}
$$
$$
= -(H_{\hat{\theta}}^{-1}g_\theta^{\boldsymbol{x}} - H_{\hat{\theta}}^{-1}g_\theta^{\boldsymbol{x}'}) \tag{35}
$$
$$
= H_{\hat{\theta}}^{-1}(g_\theta^{\boldsymbol{x}'} - g_\theta^{\boldsymbol{x}}). \tag{36}
$$

*Based on Eq.* (36)*, the significance of this gap can be estimated by*

$$
\|\Delta\mathcal{I}_{up,params}\| = \left\|H_{\hat{\theta}}^{-1}(g_\theta^{\boldsymbol{x}'} - g_\theta^{\boldsymbol{x}})\right\| \tag{37}
$$
$$
\leq \left\|H_{\hat{\theta}}^{-1}\right\|\left\|g_\theta^{\boldsymbol{x}'} - g_\theta^{\boldsymbol{x}}\right\| \tag{38}
$$
$$
= \left\|H_{\hat{\theta}}^{-1}\right\|\|\Delta g_\theta\|. \tag{39}
$$

*According to Proposition 3 and the definition of $H_{\hat{\theta}}$, we can derive that*

$$
\lim_{\epsilon\to 0}\|\Delta\mathcal{I}_{up,params}\| \leq 0. \tag{40}
$$

*Finally, $\lim_{\epsilon\to 0}\|\Delta\mathcal{I}_{up,params}\| = 0$.*

## 6. The General Workflow of Our Proposed Method

To better understand the workflow of our proposed method, we summarize the detailed algorithm in Algorithm 1. The total epoch is 200, and no warm-up schedule is used. The Adam optimizer is used with a weight decay of $1e-4$ and an initial learning rate of $3e-4$.

---

**Algorithm 1** The general workflow.

**Require:** a dataset $\mathcal{D}$ of images and label pairs $(\boldsymbol{x}, y)$, the expected selection ratio $s_r$, total epoch number $T$, classification model $f$, batch size $B$
**Ensure:** Importance scores $\mathcal{IS}$ for sample selection
1: **for** $t = 0{:}T - 1$ **do**
2:     Calculate the distance matrix $\boldsymbol{D}_k$ for each class in current feature space obtained by $f_t$;
3:     Calculate the degree of $\epsilon$-cover for each sample according to Eq. (6);
4:     Sample a mini-batch $\{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^B$ from $\mathcal{D}$;
5:     Determine the corresponding importance scores $\mathcal{IS}$ for each sample $\boldsymbol{x}_i$ using $\theta_{at}$;
6:     Calculate the reward signal $r_1$ according to Eq. (5);
7:     Calculate the reward signal $r_2$ according to Eq. (7);
8:     Calculate the general reward value $r = r_1 + r_2$;
9:     Update $\theta_a$ and $\theta_c$ according to Eq. (8) and Eq. (9), respectively;
10:     Update model $f_t$ on the current mini-batch using vanilla cross-entropy loss;
11: **end for**

---

## 7. Choice of A2C Network

Table 6. The A2C network architecture details.

|  | index | Layer | Dimension |
|---|---|---|---|
| | 1 | Linear | (512, 512) |
| Actor | 2 | Linear | (512, 256) |
| | 3 | Linear | (256, 1) |
| | 1 | Linear | (512, 512) |
| Critic | 2 | Linear | (512, 256) |
| | 3 | Linear | (256, 1) |

In Table 6, we provide the details of the A2C network. Both the actor and critic in A2C consist of three linear layers. Therefore, forwarding and updating the A2C network can be very efficient.

## 8. Complexity Analysis

According to the algorithm pipeline in Algorithm 1, the main computational costs can be divided into three components: 1) distance metric calculation, 2) degree of $\epsilon$-cover calculation, and 3) the RL network forward passes and updates. The complexities of the first two steps are $O\left(N_k^2 d\right)$ and $O\left(N_k^2\right)$, respectively, where $N_k$ is the number of samples in class $k$ and $d$ is the feature dimension (e.g., 512 for ResNet-18). Since a precise theoretical analysis of the computational complexity of the A2C algorithm is challenging due to its inherent nature, we provide some insights into its computational cost. The structure of A2C algorithm is

Table 7. The reductions in training costs with ImageNet-1k selected datasets compared to vanilla training. The reported results are the average $\pm$ std across five independent runs.

| Selection Ratio | 60% | 70% | 80% | 90% | 100% |
|---|---|---|---|---|---|
| ResNet-50 | $-35.45_{\pm0.3}$ | $-26.15_{\pm0.3}$ | $-18.23_{\pm0.3}$ | $-10.16_{\pm0.3}$ | -0 |

very simple (in Table 6): both the policy and critic networks consist of only a few linear layers, making it more computationally efficient compared to previous methods. Thus, our method can achieve competitive training efficiency compared to other baselines in Figure 3. As a result, our proposed method achieves a better trade-off between computational costs and performance, making it the best-performing method with competitive computational costs.

## 9. Implementation Details

Our RL module follows the design and parameter settings from [45], without introducing any additional hyperparameters. For experiments on CIFAR-10 and CIFAR-100, following [72, 77, 88], we train ResNet-50 models for 200 epochs with a batch size of 256 and a 0.1 learning rate with a cosine annealing learning rate decay strategy, an SGD optimizer with a momentum of 0.9, and weight decay of 5e-4. Data augmentation of random crop and random horizontal flip is added. For experiments on Tiny-ImageNet, following [72], we adopt a batch size of 256, an SGD optimizer with a momentum of 0.9, weight decay of 1e-4, and an initial learning rate of 0.1. The learning rate is divided by 10 after the 30th and the 60th epoch. The total number of epochs is 90. In each experiment, we perform three independent random trials. For experiments on ImageNet-1k, following [61, 72, 88], the VISSL library [14] is exploited. We adopt a base learning rate of 0.01, a batch size of 256, an SGD optimizer with a momentum of 0.9, and a weight decay of 1e-3. Because of the huge computational cost, the experiment in each case is performed once. Note that the methods Glister and CG-Score incur high computational and memory costs due to the iterative solving of the bi-level optimization problem [28] and the calculation of large Gram matrix inversions [49] for subset selection, respectively. Thus, they are not compared on ImageNet-1k.

## 10. Data Selection Improves Training Efficiency

Data selection substantially enhances the efficiency of model training efficiency for subsequent tasks. Specifically, once the selected datasets are obtained, only a subset needs to be stored as a replacement for the full dataset, leading to savings in memory costs. Moreover, the reduction in training data volume translates to diminished training costs. In Table 7, we show the reductions in training deep models on

ImageNet-1k. It can be seen that the reductions in training costs are proportional to the number of selected data. Meanwhile, it is important to note that with high data selection ratios, the generalization performance is nearly lossless or even enhanced, as shown in Section 5. Therefore, the selected datasets offer practical benefits.

## 11. Experiment Results on CIFAR-10

Due to space constraints, we present experiment results on CIFAR-10. As shown in Table 8, our method can achieve superior results. Notably, with relatively high selection ratios (e.g., $\geq 90\%$), our method exhibits performance surpassing that of the full dataset. Moreover, when the data selection ratio is very low, e.g., 20% and 30%, our method substantially improves over existing baselines. Therefore, through extensive experiment results alongside those obtained from other benchmark datasets as delineated in Section 5, we demonstrate the promising efficacy of our proposed method.

## 12. More Experimental Results on Vision Transformer

To further demonstrate the superiority of our proposed method, we employ Vision Transformer [10] to train with selected datasets. Following [72], the implementation is based on the public Github repository [1], where ViT small is used. We systematically evaluate our approach across different selection ratios using the CIFAR-10 dataset. Experimental results in Table 9 demonstrate the notable performance gains achieved by our method with ViT compared to other baselines.

## 13. Limitation and Future Work

First, the proposed method focuses on optimizing sample-wise importance scores for selection. This may limit its applicability to extremely large-scale datasets, which is also a primary challenge for all score-based data selection methods. Future work should emphasize extending existing data selection approaches to such large-scale datasets. Second, the proposed method drops samples that are more likely to be $\epsilon$-covered by others. This is based on an important assumption that the amount of noise in datasets is limited. However, if the amount of noise is dominant, the usefulness of our method is not guaranteed, as noisy samples are more likely to be outliers and less likely to be covered by normal samples. Therefore, it is very necessary to develop a variant that adapts to high-noise conditions with theoretical guarantees in future work. Lastly, this paper evaluates the performance of the proposed method on classification tasks. Future work should, therefore, extend its application

---

[1] https://github.com/kentaroy47/vision-transformers-cifar10

Table 8. Test accuracy (%) on CIFAR-10 with ResNet-50.

| Method / Selection ratio | 20% | 30% | 40% | 60% | 70% | 80% | 90% | 100% |
|---|---|---|---|---|---|---|---|---|
| Random | $84.12_{\pm1.53}$ | $90.34_{\pm0.39}$ | $92.71_{\pm0.38}$ | $94.43_{\pm0.37}$ | $95.02_{\pm0.29}$ | $95.55_{\pm0.14}$ | $95.89_{\pm0.11}$ | $96.12_{\pm0.12}$ |
| EL2N | $70.32_{\pm0.74}$ | $87.48_{\pm0.80}$ | $89.23_{\pm0.61}$ | $94.43_{\pm0.27}$ | $95.17_{\pm0.27}$ | $95.55_{\pm0.18}$ | $96.01_{\pm0.20}$ | $96.12_{\pm0.12}$ |
| MoSo | $83.33_{\pm0.47}$ | $89.17_{\pm0.14}$ | $92.47_{\pm0.14}$ | $94.69_{\pm0.20}$ | $95.50_{\pm0.00}$ | $95.93_{\pm0.01}$ | $96.26_{\pm0.02}$ | $96.12_{\pm0.12}$ |
| GraNd | $79.23_{\pm0.84}$ | $87.88_{\pm0.90}$ | $92.17_{\pm0.73}$ | $94.14_{\pm0.47}$ | $95.19_{\pm0.12}$ | $95.35_{\pm0.38}$ | $95.96_{\pm0.05}$ | $96.12_{\pm0.12}$ |
| Glister | $79.23_{\pm0.55}$ | $87.88_{\pm0.49}$ | $92.17_{\pm0.34}$ | $95.03_{\pm0.13}$ | $95.61_{\pm0.05}$ | $95.98_{\pm0.17}$ | $96.34_{\pm0.02}$ | $96.12_{\pm0.12}$ |
| Herding | $78.42_{\pm0.78}$ | $87.77_{\pm0.66}$ | $89.40_{\pm0.54}$ | $89.12_{\pm0.35}$ | $92.11_{\pm0.13}$ | $93.92_{\pm0.36}$ | $95.50_{\pm0.13}$ | $96.12_{\pm0.12}$ |
| CG-Score | $80.50_{\pm1.23}$ | $89.35_{\pm0.87}$ | $92.73_{\pm0.37}$ | $95.19_{\pm0.23}$ | $95.87_{\pm0.17}$ | $\mathbf{95.99}_{\pm0.16}$ | $96.16_{\pm0.15}$ | $96.12_{\pm0.12}$ |
| Forgetting | $67.58_{\pm1.05}$ | $88.12_{\pm1.40}$ | $\mathbf{93.61}_{\pm0.87}$ | $95.17_{\pm0.25}$ | $95.85_{\pm0.20}$ | $95.46_{\pm0.27}$ | $95.85_{\pm0.37}$ | $96.12_{\pm0.12}$ |
| Moderate-DS | $81.75_{\pm0.38}$ | $90.94_{\pm0.27}$ | $92.79_{\pm0.31}$ | $94.69_{\pm0.24}$ | $95.26_{\pm0.30}$ | $95.73_{\pm0.19}$ | $96.17_{\pm0.15}$ | $96.12_{\pm0.12}$ |
| Self-sup. prototypes | $84.60_{\pm1.01}$ | $90.07_{\pm1.14}$ | $92.64_{\pm0.93}$ | $94.42_{\pm0.72}$ | $94.98_{\pm0.61}$ | $95.87_{\pm0.53}$ | $95.95_{\pm0.44}$ | $96.12_{\pm0.12}$ |
| Ours | $\mathbf{88.47}_{\pm0.79}$ | $\mathbf{91.27}_{\pm0.28}$ | $93.02_{\pm0.32}$ | $\mathbf{95.39}_{\pm0.29}$ | $\mathbf{96.08}_{\pm0.19}$ | $95.84_{\pm0.18}$ | $\mathbf{96.38}_{\pm0.10}$ | $96.12_{\pm0.12}$ |

Table 9. The test accuracy (%) on CIFAR-10 with ViT-small. ResNet $\rightarrow$ ViT.

| Method/Selection Ratio | 20% | 30% | 40% | 60% | 70% | 80% | 90% | 100% |
|---|---|---|---|---|---|---|---|---|
| Random | $67.98_{\pm0.29}$ | $71.99_{\pm0.12}$ | $74.69_{\pm0.26}$ | $78.98_{\pm0.28}$ | $80.30_{\pm0.36}$ | $81.33_{\pm0.10}$ | $82.63_{\pm0.18}$ | $84.00_{\pm0.32}$ |
| EL2N | $68.34_{\pm0.18}$ | $72.03_{\pm0.52}$ | $74.85_{\pm0.24}$ | $79.35_{\pm0.09}$ | $80.73_{\pm0.08}$ | $81.62_{\pm0.08}$ | $82.90_{\pm0.09}$ | $84.00_{\pm0.32}$ |
| MoSo | $67.15_{\pm0.19}$ | $71.80_{\pm0.18}$ | $74.88_{\pm0.19}$ | $79.45_{\pm0.11}$ | $80.27_{\pm0.23}$ | $81.82_{\pm0.15}$ | $82.92_{\pm0.34}$ | $84.00_{\pm0.32}$ |
| GraNd | $67.74_{\pm0.25}$ | $71.99_{\pm0.32}$ | $75.24_{\pm0.12}$ | $79.22_{\pm0.06}$ | $80.59_{\pm0.19}$ | $81.53_{\pm0.18}$ | $82.72_{\pm0.08}$ | $84.00_{\pm0.32}$ |
| Glister | $61.16_{\pm0.07}$ | $67.36_{\pm0.05}$ | $71.77_{\pm0.14}$ | $78.33_{\pm0.01}$ | $79.84_{\pm0.27}$ | $81.33_{\pm0.05}$ | $82.65_{\pm0.09}$ | $84.00_{\pm0.32}$ |
| Herding | $64.97_{\pm0.66}$ | $70.18_{\pm0.13}$ | $73.27_{\pm0.19}$ | $76.08_{\pm0.19}$ | $78.53_{\pm0.45}$ | $80.31_{\pm0.01}$ | $82.08_{\pm0.02}$ | $84.00_{\pm0.32}$ |
| CG-Score | $57.21_{\pm0.11}$ | $66.38_{\pm0.09}$ | $71.89_{\pm0.07}$ | $79.09_{\pm0.29}$ | $80.96_{\pm0.05}$ | $82.02_{\pm0.23}$ | $82.91_{\pm0.05}$ | $84.00_{\pm0.32}$ |
| Forgetting | $49.50_{\pm0.14}$ | $58.83_{\pm0.53}$ | $67.43_{\pm0.18}$ | $77.87_{\pm0.32}$ | $80.86_{\pm0.08}$ | $81.90_{\pm0.31}$ | $82.69_{\pm0.20}$ | $84.00_{\pm0.32}$ |
| Moderate-DS | $68.69_{\pm0.40}$ | $72.36_{\pm0.11}$ | $75.44_{\pm0.53}$ | $79.54_{\pm0.19}$ | $81.28_{\pm0.13}$ | $81.98_{\pm0.16}$ | $82.61_{\pm0.27}$ | $84.00_{\pm0.32}$ |
| Self-sup. prototypes | $67.97_{\pm0.17}$ | $72.08_{\pm0.32}$ | $75.38_{\pm0.05}$ | $79.24_{\pm0.16}$ | $80.34_{\pm0.21}$ | $81.66_{\pm0.25}$ | $82.86_{\pm0.19}$ | $84.00_{\pm0.32}$ |
| Ours | $\mathbf{69.21}_{\pm0.21}$ | $\mathbf{73.28}_{\pm0.21}$ | $\mathbf{75.48}_{\pm0.11}$ | $\mathbf{80.52}_{\pm0.21}$ | $\mathbf{81.73}_{\pm0.19}$ | $\mathbf{82.15}_{\pm0.26}$ | $\mathbf{82.96}_{\pm0.06}$ | $84.00_{\pm0.32}$ |

to a broader range of tasks, such as fine-grained classification, semantic segmentation, and object detection, and further explore its application on multimodal data.