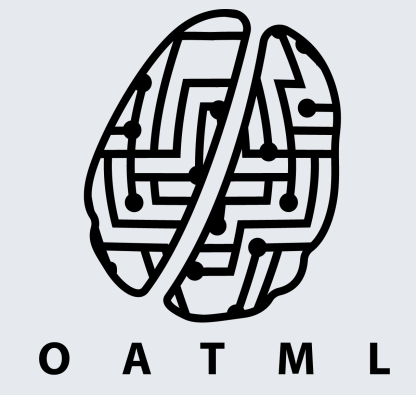


# GoldiProx: Prioritized learning on points that are learnable, worth learning, and not yet learned



1

Soren Mindermann<sup>1\*</sup>, Muhammed Razzak<sup>1\*</sup>, Winnie Xu<sup>2, 3\*</sup>

Andreas Kirsch<sup>1</sup>, Mrinank Sharma<sup>1</sup>, Adrien Morisot<sup>3</sup>, Aidan N. Gomez<sup>1,2,3</sup>, Sebastian Farquhar<sup>1</sup>, Jan Brauner<sup>1</sup>, Yarin Gal<sup>1</sup>



2



3

## Contributions

**We developed an information theoretic approach to efficient model training in the large data regime.**

Maximal information gain about the labels of a validation set:

- *Pointwise predictive information gain* as acquisition function.
- *GoldiProx Selection*: online batch selection via a proxy.

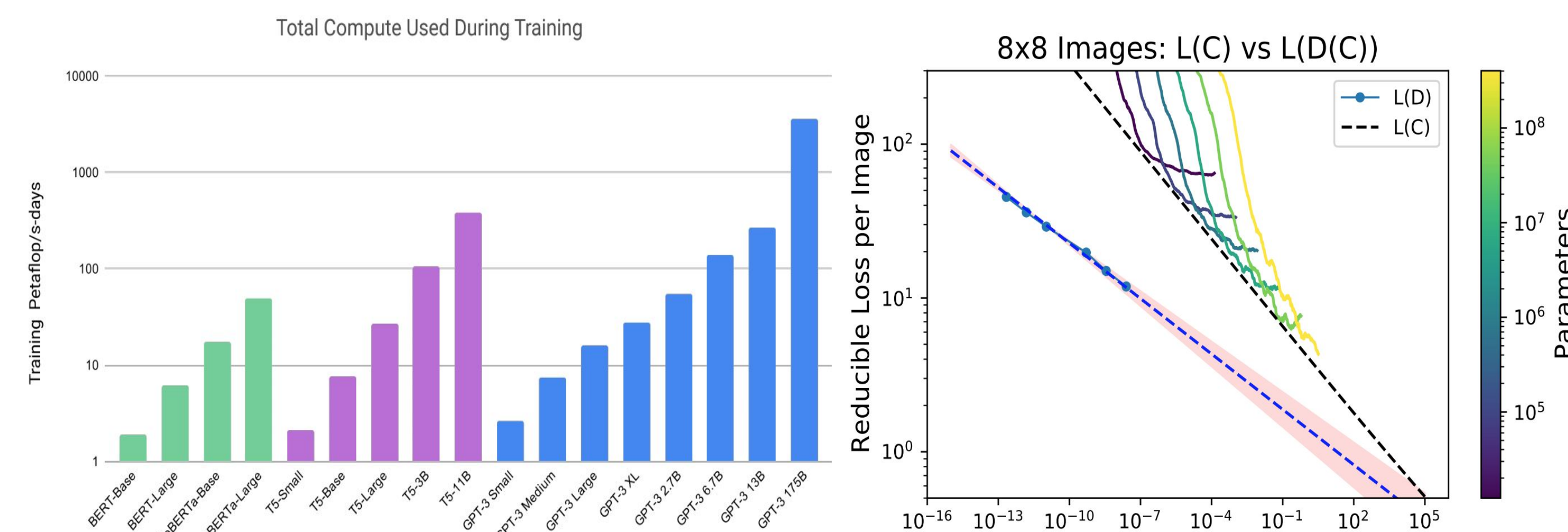
The *reducible held-out cross-entropy loss* as an easily implementable and novel objective function.

Prioritizing learnable, but informative data relevant to the eval task.

## A Transition from Mini-Batch to Mini-Epoch

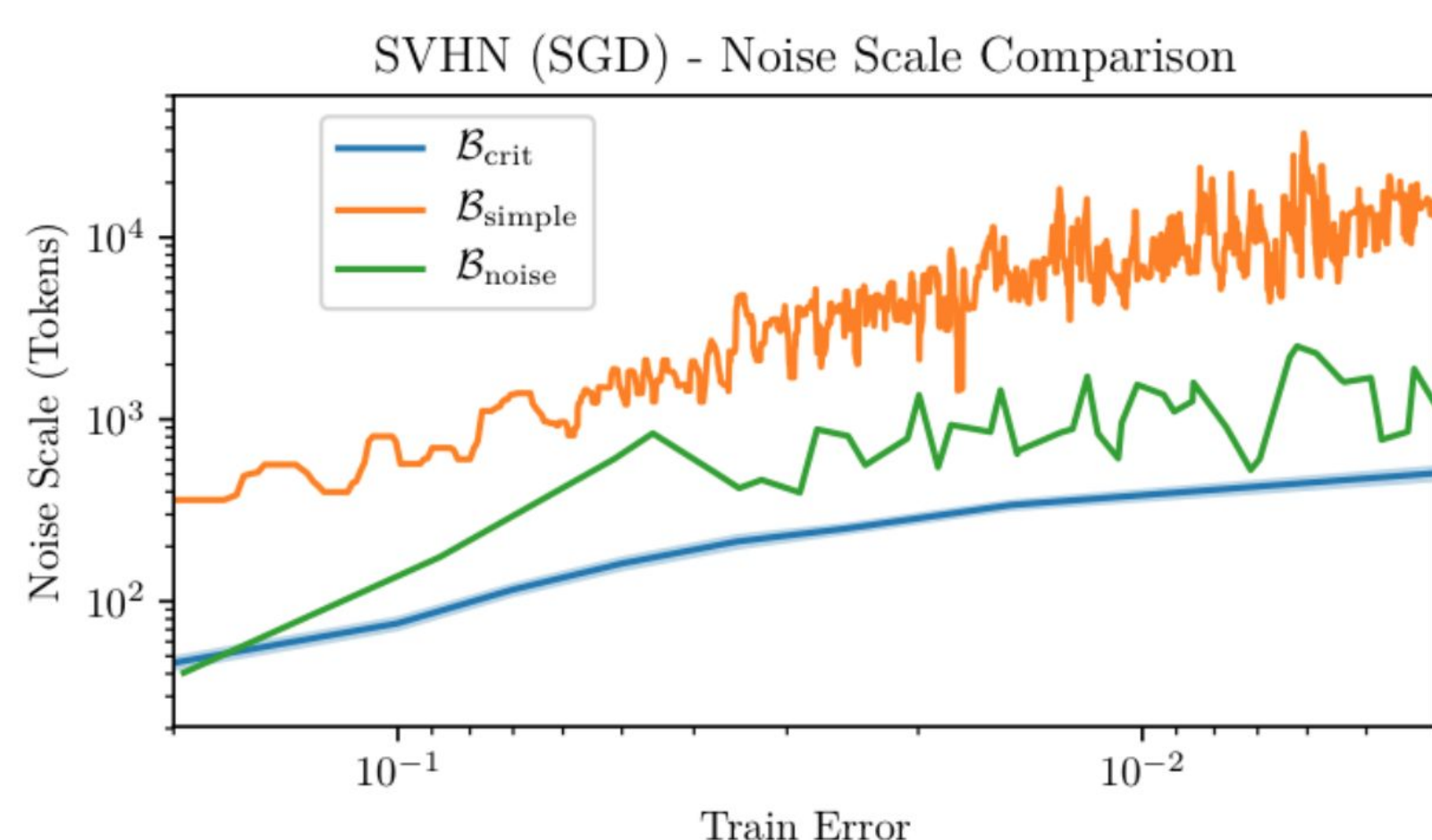
**Training is expensive and most data is uninformative:**

1. Redundant: already learned, “easy” points
2. Ambiguous: non-learnable, “noisy” and unpredictable labels
3. Atypical: not worth learning (or prioritizing), “uncertain” labels



**Common methods of data selection have shortfalls:**

1. Online batch selection: train only on “hard” points by sampling without replacement within in epoch to select a training batch  
→ Skips redundant points, but expensive forward pass
2. Curriculum learning: prioritize “easy” points with low noise before training on all points with equal weight  
→ Better convergence but fails to de-prioritize redundant points
3. Active learning: selects points where model is uncertain of label  
→ Selects informative, non-noisy, but test-time irrelevant points



Select cheaply with a small proxy, train larger model on filtered sequence

## Label-Aware Training Objective

**Expected Predictive Information Gain (EPIG)**: select informative unlabeled points by computing the info gain / mutual information between  $Y^{\text{val}}$  and  $Y$  since labels  $y$  for  $x$  are unavailable:

$$I[Y^{\text{val}}; Y | \mathbf{x}^{\text{val}}, x, \mathcal{D}_t] = H[Y^{\text{val}} | \mathbf{x}^{\text{val}}, \mathcal{D}_t] - H[Y^{\text{val}} | \mathbf{x}^{\text{val}}, \mathcal{D}_t, x, Y]$$

**Pointwise Predictive Information Gain** via **pointwise mutual information (pmi)**: remove the need to compute expectations over  $Y$  and  $Y^{\text{val}}$  since the labels  $y^{\text{val}}$  and  $y$  are available, i.e. a label-aware mutual info metric. Difference of pointwise conditional entropies (PCE):

$$\text{pmi}[y^{\text{val}}; y | \mathbf{x}^{\text{val}}, x, \mathcal{D}_t] = \underbrace{h[y^{\text{val}} | \mathbf{x}^{\text{val}}, \mathcal{D}_t]}_{\text{PCE prior to } (x, y)} - \underbrace{h[y^{\text{val}} | \mathbf{x}^{\text{val}}, \mathcal{D}_t, x, y]}_{\text{PCE after seeing } (x, y)}$$

## Reducible Held-Out Cross-Entropy Loss

**Measures the uncertainty reduction about the validation labels  $y^{\text{val}}$  due to observing a point in training  $(x, y)$ ,**

$$\underbrace{h[y | x, \mathcal{D}_t]}_{\text{untrained model cross-entropy}} - \underbrace{h[y | x, \mathbf{x}^{\text{val}}, y^{\text{val}}]}_{\text{irreducible cross entropy (loss of proxy and validation set, i.e. } D_t \ll D_{\text{val}})}$$

so that a filtered dataset  $D_t$  represents a viable sequence of training examples that can be used to train a new model while also maximizing the validation likelihood of the pre-trained proxy model.

## GoldiProx Selection: Improved Training Efficiency

1: **Input**: Initial parameters  $\theta_{\text{small}}^0$  and  $\theta_{\text{big}}^0$ , learning rate  $\eta$ , small model  $p(y | x, \theta_{\text{val}})$  trained on validation set, batch size  $|b|$ , large batch size  $|B| > |b|$ .

**Train a smaller model normally with negative log likelihood objective**

2: **for**  $i$  in training set **do**  
3: IrreducibleLoss[i]  $\leftarrow L(y_i, p(y_i | x_i, \theta_{\text{val}}))$   
4: **end for**  
5: Sequence  $\leftarrow []$

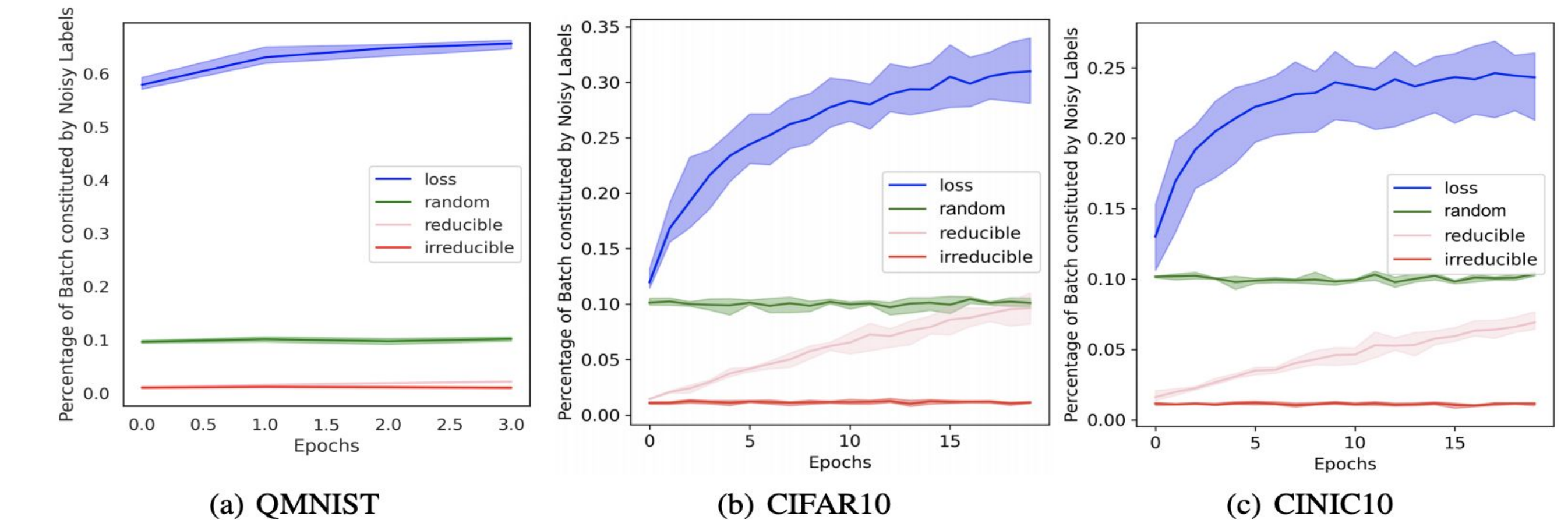
**Take a gradient step on a mini-batch, based on reducible loss**

16: **for**  $t = 1, 2, \dots$  **do**  $\triangleright$  Train large model  
17: Load  $b_t \leftarrow \text{Sequence}[t]$   
18:  $g_t \leftarrow$  mini-batch gradient on  $b_t$  for  $\theta_{\text{big}}^t$   
19:  $\theta_{\text{big}}^{t+1} \leftarrow \theta_{\text{big}}^t - \eta g_t$   
20: **end for**

6: **for**  $t = 1, 2, \dots$  **do**  $\triangleright$  Select data with small model  
7: Randomly select a large batch  $B_t$  of size  $|B|$ .  
8:  $\forall i \in B_t$ , compute Loss[i], the loss of point  $i$  given parameters  $\theta_{\text{small}}^t$   
9:  $\forall i \in B_t$ , compute ReducibleLoss[i]  $\leftarrow \text{Loss}[i] - \text{IrreducibleLoss}[i]$   
10:  $b_t \leftarrow$  top- $|b|$  samples in  $B_t$  in terms of ReducibleLoss.  
11:  $g_t \leftarrow$  mini-batch gradient on  $b_t$  for  $\theta_{\text{small}}^t$   
12:  $\theta_{\text{small}}^{t+1} \leftarrow \theta_{\text{small}}^t - \eta g_t$   
13: Append  $b_t$  to Sequence.  
14: **end for**

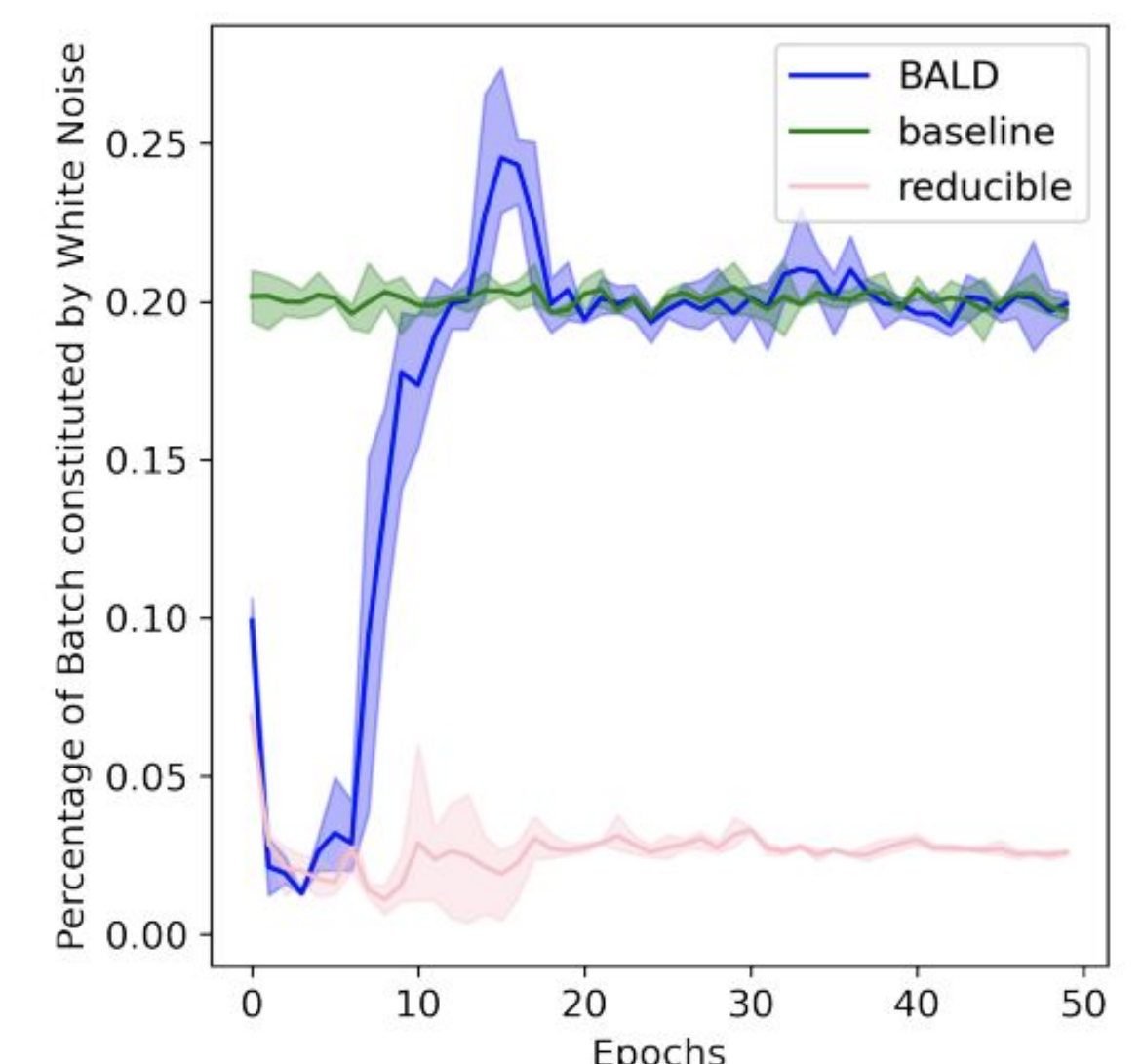
## Performant Models with Less Data and Iterations

High loss points are the points mislabelled by their noise distribution



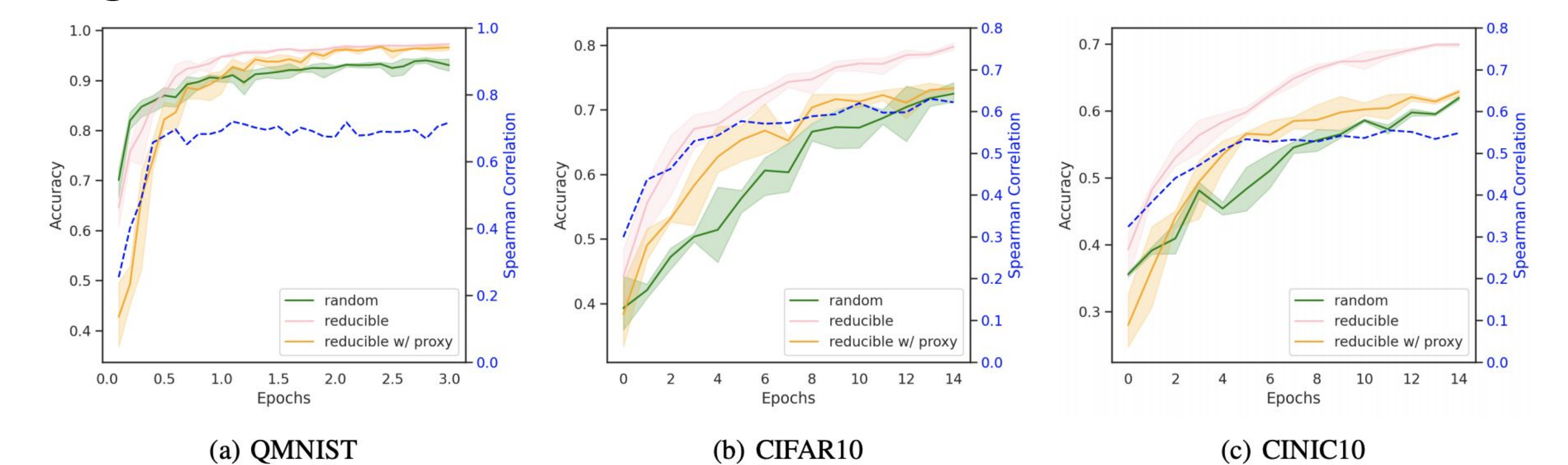
## Reduced Uncertainty for the Prediction Task

Bayesian Active Learning with Disagreement (BALD) selects 20% of noisy samples perhaps since information can be gained from low support regions of input space. This is problematic in high kurtosis input distributions, as data that is unlikely to appear at test time is not worth learning. Reducible loss triumphs.



## High Correlation between Proxy and Reducible Model

**Effective transfer of reducible loss training with small proxy to larger models of different architectures. See rank correlations.**



Model	Num of Params.	FLOPS	Model	Num of Params.	FLOPS
3-layer MLP with 128 hidden units	135 k	136 k	Small CNN	0.538 M	0.019 G
3-layer MLP with 512 hidden units	932 k	935 k	ResNet-18	11.17 M	0.557 G

This implies that a point learned by a converged model (i.e. trained on  $D$ , the original dataset) may still be informative to a model trained on  $D_t$ , the prioritized sequence of examples.

## References

- [1] Kirsch, A., van Amersfoort, J., & Gal, Y. (2019). BatchBALD: Efficient and Diverse Batch Acquisition for Deep Bayesian Active Learning. arXiv.org.
- [2] Loshchilov, Ilya, and Frank Hutter. "Online batch selection for faster training of neural networks." arXiv preprint arXiv:1511.06343 (2015).
- [3] Houlisby, Neil, et al. "Bayesian active learning for classification and preference learning." arXiv preprint arXiv:1112.5745 (2011).
- [4] Kirsch, Andreas, Tom Rainforth, and Yarin Gal. "Active Learning under Pool Set Distribution Shift and Noisy Data." arXiv preprint arXiv:2106.11719 (2021).