

A DERIVATIONS

Here, we show the equivalence of Eq. (3) and Eq. (4) for i.i.d. data:

$$h[\mathbf{y}^{\text{ho}} | \mathbf{x}^{\text{ho}}, \mathcal{D}_t] - h[\mathbf{y}^{\text{ho}} | \mathbf{x}^{\text{ho}}, \mathcal{D}_t \cup (x, y)] \quad (7)$$

Definition

$$= -\log p(\mathbf{y}^{\text{ho}} | \mathbf{x}^{\text{ho}}, \mathcal{D}_t) + \log p(\mathbf{y}^{\text{ho}} | y, \mathbf{x}^{\text{ho}}, x, \mathcal{D}_t) \quad (8)$$

Bayes rule

$$= -\log p(\mathbf{y}^{\text{ho}} | \mathbf{x}^{\text{ho}}, \mathcal{D}_t) + \log \frac{p(y | \mathbf{y}^{\text{ho}}, \mathbf{x}^{\text{ho}}, x, \mathcal{D}_t) p(\mathbf{y}^{\text{ho}} | \mathbf{x}^{\text{ho}}, x, \mathcal{D}_t)}{p(y | \mathbf{x}^{\text{ho}}, x, \mathcal{D}_t)} \quad (9)$$

Rearrange symbol order

$$= -\log p(\mathbf{y}^{\text{ho}} | \mathbf{x}^{\text{ho}}, \mathcal{D}_t) + \log \frac{p(y | x, \mathbf{y}^{\text{ho}}, \mathbf{x}^{\text{ho}}, \mathcal{D}_t) p(\mathbf{y}^{\text{ho}} | \mathbf{x}^{\text{ho}}, x, \mathcal{D}_t)}{p(y | x, \mathbf{x}^{\text{ho}}, \mathcal{D}_t)} \quad (10)$$

Independence assumptions

$$= -\log p(\mathbf{y}^{\text{ho}} | \mathbf{x}^{\text{ho}}, \mathcal{D}_t) + \log \frac{p(y | x, \mathbf{y}^{\text{ho}}, \mathbf{x}^{\text{ho}}, \mathcal{D}_t) p(\mathbf{y}^{\text{ho}} | \mathbf{x}^{\text{ho}}, \mathcal{D}_t)}{p(y | x, \mathcal{D}_t)} \quad (11)$$

Rearrange and cancel out

$$= \log p(y | x, \mathbf{y}^{\text{ho}}, \mathbf{x}^{\text{ho}}, \mathcal{D}_t) - \log p(y | x, \mathcal{D}_t) \quad (12)$$

Definition

$$= h[y | x, \mathcal{D}_t] - h[y | x, \mathbf{x}^{\text{ho}}, \mathbf{y}^{\text{ho}}, \mathcal{D}_t]. \quad (13)$$

B RELATED WORK: EXPECTED PREDICTIVE INFORMATION GAIN

Here, we provide a brief summary of the expected predictive information gain (EPIG) $I[\mathbf{Y}^{\text{eval}}, Y | \mathbf{X}^{\text{eval}}, x], \mathcal{D}_t$ introduced by Kirsch et al. (2021) as it is relevant to our own notion of Predictive Information Content. EPIG is proposed in the context of (Bayesian) active learning, where no access to labels is assumed for the selection of informative datapoints. To emphasize the difference to our labelled held-out set, we adopt the authors' terminology of an unlabelled *evaluation set* $\{x^{\text{eval}}\}_i$. In the paper, the authors discuss the case of selecting a full batch, whereas we limit ourselves to the case of a single datapoint (x, y) for better comparability, i.e. the case of batch size $n_b = 1$.

Similar to our case, the aim is to maximize the expected information gain for these evaluation points. Given that no label access is assumed, this is approximated by the expectation over the predictions of y^{eval} and y . Furthermore, the expectation over the evaluation samples can be approximated using the joint, i.e. by

$$\mathbb{E}_{\text{P}_{\text{eval}}(x^{\text{eval}}, y^{\text{eval}})} [h[y^{\text{eval}} | x^{\text{eval}}, \mathcal{D}_t \cup (x, y)]] \quad (14)$$

$$\approx \mathbb{E}_{\text{P}_{\text{eval}}(x^{\text{eval}})} [H[Y^{\text{eval}} | x^{\text{eval}}, Y, x, \mathcal{D}_t]] \quad (15)$$

$$\geq \frac{1}{|\mathcal{D}_{\text{eval}}|} H[\mathbf{Y}^{\text{eval}} | \mathbf{x}^{\text{eval}}, Y, x, \mathcal{D}_t] \quad (16)$$

where $H[\cdot | \cdot]$ denotes the conditional entropy. It is noted that the bound is not tight in general but only in the infinite data limit when predictions become uncorrelated. That is, the authors introduce an approximation of EPIG as the expected information gain between the predictions Y for the candidate samples and the predictions for the evaluation set,

$$I[\mathbf{Y}^{\text{eval}}, Y | \mathbf{x}^{\text{eval}}, x, \mathcal{D}_t] = H[\mathbf{Y}^{\text{eval}} | \mathbf{x}^{\text{eval}}, \mathcal{D}_t] - H[\mathbf{Y}^{\text{eval}} | \mathbf{x}^{\text{eval}}, Y, x, \mathcal{D}_t] \quad (17)$$

Since the first term $H[\mathbf{Y}^{\text{eval}} | \mathbf{x}^{\text{eval}}, \mathcal{D}_t]$ is constant for a fixed \mathcal{D}_t , maximizing the approximation is equivalent to minimising the second term. Similar to our approach, it is suggested to estimate these quantity by expanding it in the opposite direction, as

$$I[\mathbf{Y}^{\text{eval}}, Y | \mathbf{x}^{\text{eval}}, x, \mathcal{D}_t] = H[Y | x, \mathcal{D}_t] - H[Y | x, \mathbf{Y}^{\text{eval}}, \mathbf{x}^{\text{eval}}, \mathcal{D}_t]. \quad (18)$$

Since it still contains an expectation over the predicted labels, this is much harder to evaluate than in our setting where access to labels is given. The authors then suggest a way to do so.

C EXPERIMENT DETAILS

For experiments on QMNIST, we use a multi-layer perceptron with 2 hidden layers and 512 units in each hidden layer. For experiments on CIFAR-10, CIFAR-100, CINIC-10 and Clothing-1M we use a ResNet-18 (He et al., 2016). All models are trained using the Adam optimizer with default PyTorch hyperparameters ($\beta_1=0.9$, $\beta_2=0.999$, and weight decay of 0.01, learning rate 0.001), a training batch size $n_b = 32$ (64 for CINIC-10) with $n_B = 320$ (640 for CINIC-10) points pre-sampled, meaning we select $\frac{n_b}{n_B} = 10\%$ of points. On all datasets except QMNIST, we train using data augmentation (random crop and horizontal flip). The irreducible loss models are trained on holdout sets (not test sets, see above). For each dataset, we select the irreducible loss model from the epoch with lowest held-out loss on \mathcal{D} (as opposed to highest accuracy). The multiple target architectures in Fig 5 were adapted from Phan (2021).

D IRREDUCIBLE HELD-OUT LOSS APPROXIMATION

To empirically justify the approximation made in Eq. 5, we train an irreducible loss model with limited capacity to improve when updated on \mathcal{D}_t . We make use of MNIST-8M (Loosli et al., 2007) and use an MLP with 2 hidden layers and 16 hidden units. MNIST-8M applies pseudo-random deformations to MNIST to increase the dataset to over 8 million images. We use the original MNIST training set to train the irreducible loss model and MNIST-8M training set, excluding the original MNIST training set, as \mathcal{D}_t . Note that the holdout set is thus much smaller than \mathcal{D}_t .

The accuracy of the irreducible loss model, with constrained capacity, remains at around 94% despite being updated on \mathcal{D}_t , and the updating has no discernible effect on the accuracy of the model trained exclusively on \mathcal{D}_t .

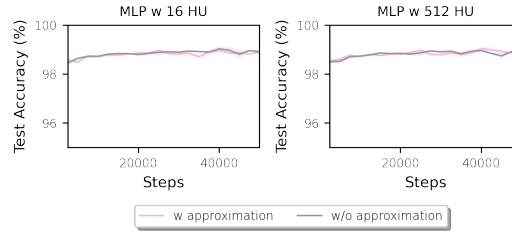


Figure 6: Reducible loss selection with and without the approximation in Eq. 5.

E SMALL IRREDUCIBLE LOSS MODELS AVOID NOISY POINTS

Here, we show one reason why selection with a small irreducible loss model still accelerates training: it still avoids points with noisy labels. Figure 7 is identical to Figure 1B, but includes selection using an irreducible loss model that uses 21x fewer parameters (and 29x fewer FLOPs) than the target model. It also has a simpler architecture (vanilla CNN versus ResNet-18). Across our datasets, selection with the smaller irreducible loss model avoids points whose labels are corrupted with uniform noise.

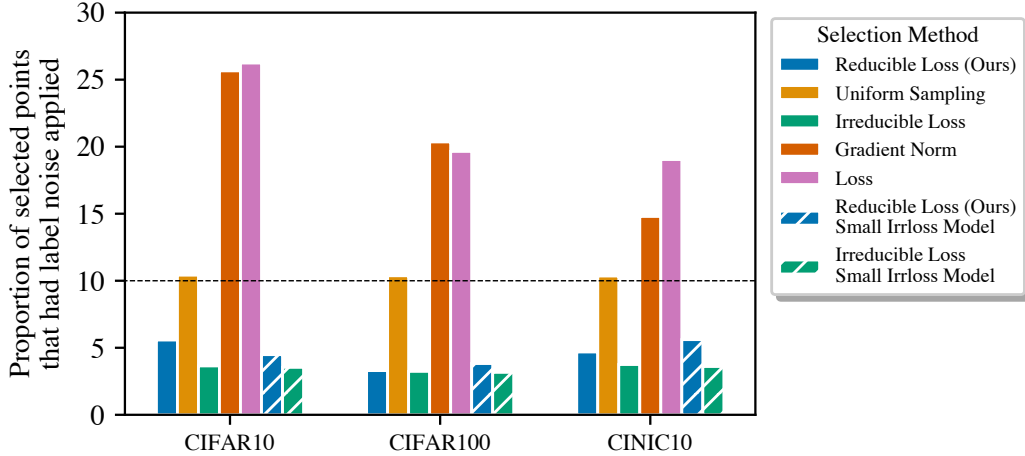


Figure 7: Proportion of selected points with label noise applied for different selection method. We added 10% uniform label noise, i.e., we randomly switched each point’s label with 10% probability. Selecting the 10% “hardest” points (high loss or gradient norm) gives batches where $> 10\%$ of the labels are mislabeled. We can avoid noisy points by selecting points that are “easy” (low irreducible loss) or have high reducible loss. This works even when approximating the irreducible loss with a model that uses 29x fewer FLOPs and has 21x fewer parameters than the target model to be trained. Results are averaged over 15 epochs.

F ABLATION OF PERCENTAGE SELECTED

Our method has a hyperparameter, the percentage $\frac{n_b}{n_B}$ of evaluated points which are selected for training. In the experiments above, this parameter was set to 0.1. We have not tuned this parameter, as we aim to analyse how well our method works “out of the box”. In fact, on some datasets, performance further improves with other values of this parameter settings of this hyperparameters seem to work better to give the highest speed. Adjusting this percentage should allow practitioners to specify their preferred tradeoff between training time and computation, where a low percentage corresponds to a lower training time. Training further accelerates when we lower the percentage from 10% to 5% and decelerates when we increase it to 20% or 50%. For these experiments, we keep $n_b = 32$ and adapt n_B accordingly. The percentage $\frac{n_b}{n_B}$ of datapoints selected per batch has different effects across datasets as shown in Figure 8.

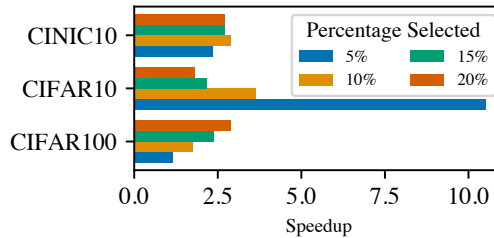


Figure 8: Varying the percent of data points selected in each training batch. Average over 3 random seeds.

G ACTIVE LEARNING BASELINES

We compare our method to typical methods used in the Active Learning (AL) literature. Note that our method is label-aware, while active learning acquires datapoints without using labelled information. We consider the following baselines, which select the top- k points using an acquisition function, $\alpha(x)$:

- Bayesian Active Learning by Disagreement (Houlsby et al., 2011) with $\alpha(x) = H[y | x, \mathcal{D}_t] - \mathbb{E}_{p(\theta|\mathcal{D}_t)} [H[y | x, \theta]]$.
- (Average) conditional entropy, $\alpha(x) = \mathbb{E}_{p(\theta|\mathcal{D}_t)} [H[y | x, \theta]]$, where the average is taken over the model parameter posterior.
- (Average) entropy, $\alpha(x) = H[y | x, \mathcal{D}_t]$.
- Loss minus conditional entropy $\alpha(x) = h[y | x, \theta] - \mathbb{E}_{p(\theta|\mathcal{D}_t)} [H[y | x, \theta]]$. This uses the (average) conditional entropy as an estimate of how noisy datapoint x is—points with high noise are deprioritised.

We additionally compare our method to uniform sampling. We run all baselines on MNIST and CIFAR10. Note that several of these active learning baselines consider epistemic uncertainty; that is, uncertainty in predictions driven by uncertainty in the model parameters. This mandates performing (approximate) Bayesian inference. We use Monte-Carlo Dropout (Gal & Ghahramani, 2016) to perform approximate inference. For MNIST, we use a 2 hidden layer MLP with 512 hidden units per hidden layer, and a dropout probability of a 0.5. For experiments on CIFAR10, we use a small-scale CNN with dropout probability 0.05 (the dropout probability follows (Osawa et al., 2019)).

Figure 9 shows training curves for our method, uniform sampling, and the active learning baselines. Our method accelerates training across both datasets. The active learning methods accelerate training for MNIST but not for CIFAR10. This highlights that active learning methods, if naively applied to online batch selection, may not accelerate model training.

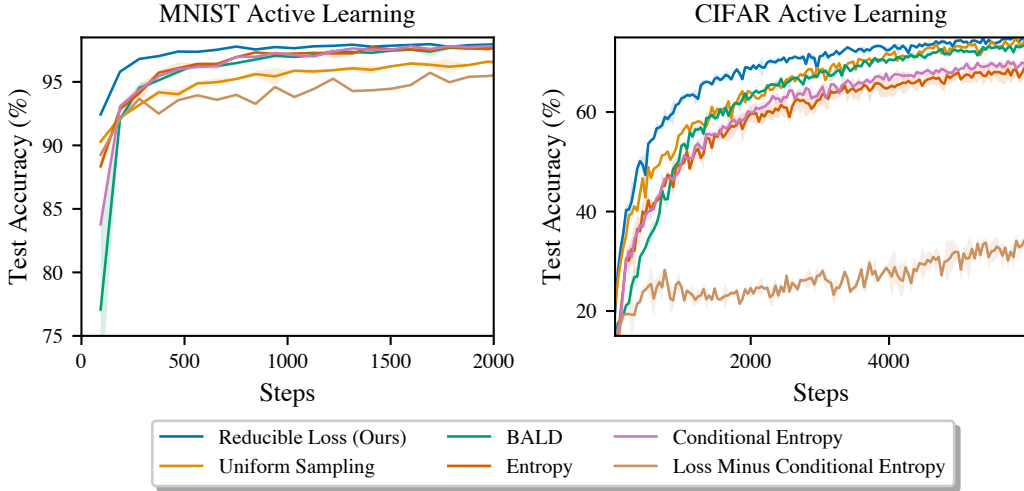


Figure 9: Training curves for several active learning baselines on the MNIST and CIFAR10 datasets.