# Stochastic Batch Acquisition:
# A Simple Baseline for Deep Active Learning

**Anonymous authors**
**Paper under double-blind review**

## Abstract

We examine a simple stochastic strategy for adapting well-known single-point acquisition functions to allow batch active learning. Unlike acquiring the top-$K$ points from the pool set, score- or rank-based sampling takes into account that acquisition scores change as new data are acquired. This simple strategy for adapting standard single-sample acquisition strategies performs just as well as compute-intensive state-of-the-art batch acquisition functions, like BatchBALD or BADGE while using orders of magnitude less compute. In addition to providing a practical option for machine learning practitioners, the surprising success of the proposed method in a wide range of experimental settings raises a difficult question for the field: when are these expensive batch acquisition methods pulling their weight?

## 1 Introduction

Active learning is a widely used strategy for efficient learning in settings where unlabelled data are plentiful, but labels are expensive (Atlas et al., 1989; Settles, 2010). For example, labels for medical image data may require highly trained annotators, and when labels are the results of scientific experiments, each one can require months of work. Active learning uses information about unlabelled data and the current state of the model to acquire labels for those samples that are most likely to be informative.

While many acquisition schemes are designed to acquire labels one at a time (Houlsby et al., 2011; Gal et al., 2017), recent work has highlighted the importance of *batch acquisition* (Kirsch et al., 2019; Ash et al., 2020). Acquiring in a batch lets us parallelise labelling. For example, we could hire hundreds of annotators to work in parallel or run more than one experiment at once. Batch acquisition also saves compute as single-point selection also incurs the cost of retraining the model for every new data point.

Unfortunately, existing batch acquisition schemes are computationally expensive (Table 1). Intuitively, this is because batch acquisition schemes face combinatorial complexity when accounting for the interactions between possible acquisition points. Recent works (Ash et al., 2020; 2021) trade off a principled motivation with various approximations to remain tractable. A commonly used, though extreme, heuristic is to take the top-$K$ highest scoring points from an acquisition scheme designed to select a single point.

This paper examines a simple baseline for batch active learning that is competitive with methods that cost orders of magnitude more across a wide range of experimental contexts. This method is motivated by noticing that single-acquisition score methods such as BALD (Houlsby et al., 2011) act as a noisy proxy for future acquisition scores as we motivate in Figure 1. This observation leads us to stochastically acquire points following a distribution determined by the single-acquisition scores. This simple approach matches a prior state of the art for batch acquisition (BatchBALD, Kirsch et al. (2019)) despite being very simple. Indeed, this acquisition scheme has a time complexity of only $\mathcal{O}(M \log K)$ in the pool size $M$ and acquisition size $K$, just like top-$K$ acquisition.

We show empirically that the presented stochastic strategy performs as well or better than top-$K$ acquisition with almost identical computational cost on several commonly used acquisition scores, making it a strictly-better batch strategy. Strikingly, the empirical comparisons between this stochastic strategy and the evaluated more complex methods cast doubt on whether they function as well as claimed. Concretely, in this paper we:

Table 1: *Acquisition runtime (in seconds, 5 trials, $\pm$ s.d.).* The examined stochastic acquisition methods are as fast as top-$K$, and **orders of magnitude** faster than BADGE or BatchBALD. Synthetic pool set with $M = 10,000$ pool points with 10 classes. BatchBALD and BALD with 20 parameter samples.

| $K$ | Top-$K$ | **Stochastic** | BatchBALD | BADGE |
|---|---|---|---|---|
| 10 | $0.2 \pm 0.0$ | $0.2 \pm 0.0$ | $566.0 \pm 17.4$ | $9.2 \pm 0.3$ |
| 100 | $0.2 \pm 0.0$ | $0.2 \pm 0.0$ | $5,363.6 \pm 95.4$ | $82.1 \pm 2.5$ |
| 500 | $0.2 \pm 0.0$ | $0.2 \pm 0.0$ | $29,984.1 \pm 598.7$ | $409.3 \pm 3.7$ |

- examine a family of three computationally cheap stochastic batch acquisition strategies;

- demonstrate that these strategies are preferable to the commonly used top-$K$ acquisition heuristic; and

- identify the failure of two existing batch acquisition strategies to outperform this vastly cheaper and more heuristic strategy.

**Outline.** In §2, we present active learning notation and commonly used acquisition functions. We propose stochastic extensions in §3, relate them to previous work in §4, and validate them empirically in §5 on various datasets, showing that these extensions are competitive with some much more complex active learning apporaches despite being orders of magnitude computationally cheaper. Finally, we validate some of the underlying theoretical motivation in §6 and discuss limitations in §7.

## 2 Problem Setting

The stochastic approach we examine applies to batch acquisition for active learning in a pool-based setting (Settles, 2010) where we have access to a large unlabelled *pool* set, but we can only label a small subset of the points. The challenge of active learning is to use what we already know to pick which points to label in the most efficient way. Generally, we want to avoid labelling points similar to those already labelled.

**Notation.** Following Farquhar et al. (2021), we formulate active learning over *indices* instead over datapoints. This simplifies the notation. The large, initially fully unlabelled, pool set containing $M$ input points is

$$\mathcal{D}^{\text{pool}} = \{x_i\}_{i \in \mathcal{I}^{\text{pool}}}, \tag{1}$$

where $\mathcal{I}^{\text{pool}} = \{1, \dots, M\}$ is the initial full index set. We initialise a training dataset with $N_0$ randomly selected points from $\mathcal{D}^{\text{pool}}$ by acquiring their labels, $y_i$,

$$\mathcal{D}^{\text{train}} = \{(x_i, y_i)\}_{i \in \mathcal{I}^{\text{train}}}, \tag{2}$$

where $\mathcal{I}^{\text{train}}$ is the index set of $\mathcal{D}^{\text{train}}$, *initially* containing $N_0$ indices between 1 and $M$. A model of the predictive distribution, $\text{p}(y \mid x)$, can then be trained on $\mathcal{D}^{\text{train}}$.

**Active Learning.** At each acquisition step, we select additional points for which to acquire labels. Although many methods acquire one point at a time (Houlsby et al., 2011; Gal et al., 2017), one can alternatively acquire a whole batch of $K$ examples. An acquisition function $a$ takes $\mathcal{I}^{\text{train}}$ and $\mathcal{I}^{\text{pool}}$ and returns $K$ indices from $\mathcal{I}^{\text{pool}}$ to be added to $\mathcal{I}^{\text{train}}$. We then label those $K$ datapoints and add them to $\mathcal{I}^{\text{train}}$ while making them unavailable from the pool set. That is,

$$\mathcal{I}^{\text{train}} \leftarrow \mathcal{I}^{\text{train}} \cup a(\mathcal{I}^{\text{train}}, \mathcal{I}^{\text{pool}}), \tag{3}$$

$$\mathcal{I}^{\text{pool}} \leftarrow \mathcal{I}^{\text{pool}} \setminus \mathcal{I}^{\text{train}}. \tag{4}$$

A common way to construct the acquisition function is to define some scoring function, $s$, and then select the point(s) that score the highest.

**Probabilistic Model.** We assume classification with inputs $X$, labels $Y$, and a discriminative classifier $\mathrm{p}(y \mid x)$. In the case of Bayesian models, we further assume a subjective probability distribution over the parameters, $\mathrm{p}(\omega)$, and we have $\mathrm{p}(y \mid x) = \mathbb{E}_{\mathrm{p}(\omega)}[\mathrm{p}(y \mid x, \omega)]$.

**BALD.** One popular scoring function is *Bayesian Active Learning by Disagreement* (Houlsby et al., 2011) which computes the expected information gain between the predictive distribution and the parameter distribution $\mathrm{p}(\omega \mid \mathcal{D}^{\mathrm{train}})$ for a Bayesian model. For each candidate pool index, $i$, with mutual information, I, and entropy, H, the score is

$$
\begin{aligned}
s_{\mathrm{BALD}}(i; \mathcal{I}^{\mathrm{train}}) &\coloneqq \mathrm{I}[Y; \Omega \mid X = x_i, \mathcal{D}^{\mathrm{train}}] \\
&= \mathrm{H}[Y \mid X = x_i, \mathcal{D}^{\mathrm{train}}] - \mathbb{E}_{\mathrm{p}(\omega \mid \mathcal{D}^{\mathrm{train}})}[\mathrm{H}[Y \mid X = x_i, \omega]].
\end{aligned}
\tag{5}
$$

**Entropy.** Another popular scoring function is the *(predictive) entropy* (Gal et al., 2017). It does not require Bayesian models, unlike BALD, and performs worse for data with high observation noise Mukhoti et al. (2021). It is identical to the first term of the BALD score

$$
s_{\mathrm{entropy}}(i; \mathcal{I}^{\mathrm{train}}) \coloneqq \mathrm{H}[Y \mid X = x_i, \mathcal{D}^{\mathrm{train}}].
\tag{6}
$$

**Acquisition Functions.** These scoring functions were introduced for single-point acquisition:

$$
a_s(\mathcal{I}^{\mathrm{train}}) \coloneqq \underset{i \in \mathcal{I}^{\mathrm{pool}}}{\arg\max}\, s(i; \mathcal{I}^{\mathrm{train}}).
\tag{7}
$$

For deep learning in particular, single-point acquisition is computationally expensive due to retraining the model for every acquired sample. Moreover, it also means that labelling can only happen sequentially instead of in bulk. Thus, single-point acquisition functions were expanded to multi-point acquisition via acquisition batches in batch active learning. The most naive batch acquisition function selects the highest $K$ scoring points

$$
a_s^{\mathrm{batch}}(\mathcal{I}^{\mathrm{train}}; K) \coloneqq \underset{I \subseteq \mathcal{I}^{\mathrm{pool}}, |I| = K}{\arg\max} \sum_{i \in I} s(i; \mathcal{I}^{\mathrm{train}}).
\tag{8}
$$

Maximizing this sum is equivalent to taking the top-k scoring points, which cannot account for the interactions between points in an acquisition batch because individual points are scored independently. For example, if the most informative point is duplicated in the pool set, all instances will be acquired, which is likely wasteful when we assume no label noise (see also Figure 1 in Kirsch et al. (2019)). Some acquisition functions are explicitly designed for batch acquisition (Kirsch et al., 2019; Ash et al., 2020). They try to account for the interaction between points, which can improve performance relative to simply selecting the top-$K$ scoring points. However, existing methods can be computationally expensive. For example, BatchBALD rarely scales to acquisition sizes of more than 5–10 points due to its long runtime (Kirsch et al., 2019), as we evidence in Table 1.

**BatchBALD.** Kirsch et al. (2019) extend BALD to batch acquisition using the mutual information between the parameter distribution and the *joint* distribution of the predictions of multiple point in an acquistion batch: this mutual information is the expected information gain for an acquistion batch. Kirsch et al. (2019) greedily construct an acquisition batch by iteratively selecting the next unlabelled pool point that maximizes the joint score with the already selected points. This is $1 - \nicefrac{1}{e}$-optimal as the expected information gain is submodular (Krause & Golovin, 2014). They note that their approach is computationally expensive, and they only consider acquisition batches of up to size 10.

**BADGE.** Ash et al. (2020) propose *Batch Active learning by Diverse Gradient Embeddings*: it motivates its batch selection approach using a k-Determinantal Point Process (Kulesza & Taskar, 2011) based on the (inner product) similarity matrix of the scores (gradients of the log loss) using hard pseudo-labels (the highest probability class according to the model's prediction) for each pool sample. See also Kirsch & Gal (2022) for a more detailed analysis. In practice, they use the intialization step of k-MEANS++ with Euclidian distances between the scores to select an acquisition batch. BADGE is also computationally expensive.

Table 2: *Summary of stochastic acquisition variants.* Perturbing the scores $s_i$ themselves with $\epsilon_i \sim$ Gumbel$(0; \beta^{-1})$ i.i.d. yields a softmax distribution. Log-scores result in a power distribution, with assumptions that are reasonable for active learning. Using the score-ranking, $r_i$ finally is a robustifying assumption. $\beta$ is included for completeness; we use $\beta := 1$ in our experiments—except for the ablation in §6.1.

| Perturbation | Distribution | Probability mass |
|---|---|---|
| $s_i + \epsilon_i$ | Softmax | $\propto \exp \beta s_i$ |
| $\log s_i + \epsilon_i$ | Power | $\propto s_i^{\beta}$ |
| $-\log r_i + \epsilon_i$ | Soft-rank | $\propto r_i^{-\beta}$ |



Figure 1: *Acquisition scores at individual acquisition step t are only a loose proxy for later scores at $t + n$ (here: $t = 0$).* Specifically, the Spearman rank-correlation between acquisition scores on the zeroth and $n$'th time-step falls with $n$. While top-$K$ acquisition incorrectly implicitly assumes the rank-correlation remains 1, stochastic acquisitions do not. Using Monte-Carlo Dropout BNN trained on MNIST at initial 20 points and 73% initial accuracy; score ranks computed over test set.
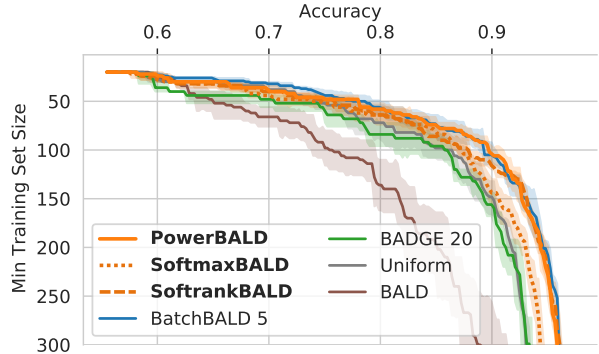
Figure 2: *Performance on Repeated-MNIST with 4 repetitions (5 trials).* **Up and to the right is better (↗).** PowerBALD outperforms (top-$K$) BALD and BADGE and is on par with BatchBALD. This is despite being orders of magnitude faster. Acquisition sizes: BatchBALD–5, BADGE–20, others–10. See Figure 10 in the appendix for an ablation study of BADGE's acquisition size.

## 3 Method

Selecting the top-$K$ points at acquisition step $t$ amounts to the assumption that the informativeness of these points is independent of each other. This leads to the pathology that if the most informative pool point is duplicated in the pool set, each instance would be selected (up to the acquisition batch size). This is clearly wrong.

Another way to view this is to split batch acquisition into multiple steps of size 1. We select the top pool sample by acquisition score and retrain the model once for each possible class label for this point. We then compute the averaged acquisition scores on the pool set given each of these models weighted by the original model's probability of each class label. We select the top pool sample by this new (averaged) score, and repeat the process, exponentially branching out as necessary. This is equivalent to the top-$K$ acquisition batch selection in BatchBALD (Kirsch et al., 2019), but it is clear that the informativeness of the pool samples will not be independent of each other: Of course, the acquisition scores for models trained with these additional points will be different from the first set of scores. After all, the purpose of active learning is to add the *most informative* points—those that will update the model the most. Yet selecting a top-$K$ batch in one step implicitly assumes that the score ranking will not change due to these points. We provide empirical confirmation in Figure 1 that, in fact, the ranking of acquisition scores at step $t$ and $t + K$ is decreasingly correlated as $K$ grows when we retrain the model for each acquired point. Figure 3 also illustrates this on MNIST. Moreover, as we will see in §6, this effect is the strongest for the most informative points.
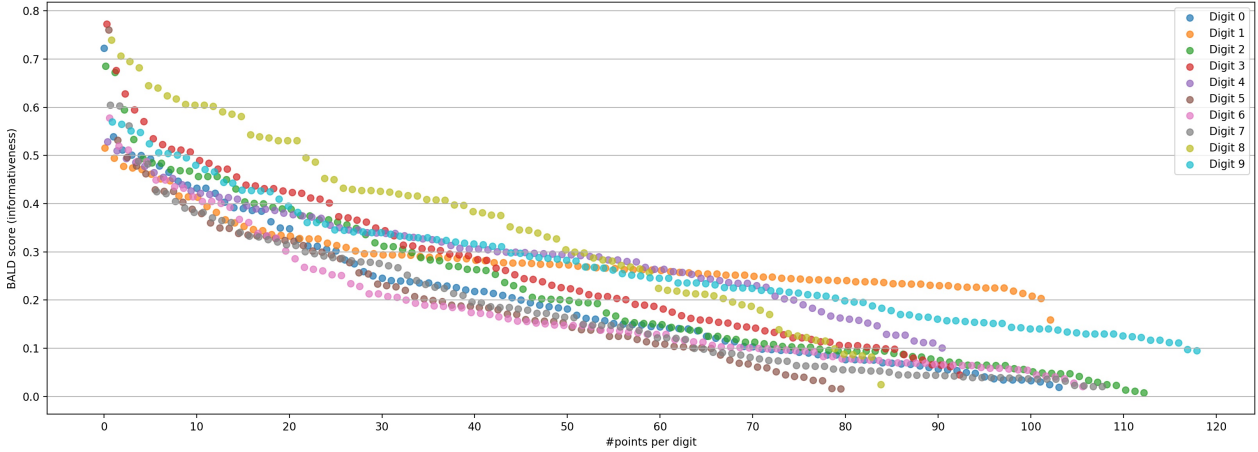
Figure 3: *BALD scores for 1000 randomly-chosen points from the MNIST dataset (hand-written digits).* The points are colour-coded by digit label and sorted by score. The model used for scoring has been trained to 90% accuracy first. If we were to pick the top scoring points (e.g. scores above 0.6), most of them would be 8s, even though we can assume that after acquiring the first couple of them our model would consider them less informative than other available data. Points are slightly jittered on the x-axis by digit label to avoid overlaps.

Instead, we investigate the use of stochastic sampling as an alternative to top-$K$ acquisition, which implicitly acknowledges the uncertainty within the batch acquisition step using a simple noise process model governing how scores change. We motivate and investigate the theory behind this in §6, but given how simple the examined methods are, this theory only obstructs their simplicity. Specifically, we examine three simple stochastic extensions of single-sample scoring functions $s(i; \mathcal{I}^{\text{train}})$ that make slightly different assumptions. These methods are compatible with conventional active learning frameworks that typically take the top-$K$ highest scoring samples. For example, it is straightforward to adapt entropy, BALD, and other scoring functions for use with these extensions.

These stochastic acquisition distributions assume that future scores differ from the current score by a perturbation. We model the noise distribution of this perturbation as the addition of Gumbel-distributed noise $\epsilon_i \sim \text{Gumbel}(0; 1)$, which is used frequently for modelling extrema.

Why is this sensible? Intuitively, to select the $k$-th point, we want to take into account how much additional information (increase in acquisition scores) the still-to-be-selected additional $K - k$ points will provide. As such we want to model the maximum over all possible additional candidate points. The maximum of a set of random variables that follow an exponential distribution is known to be well approximated by a Gumbel distribution in the sample limit (Gumbel, 1954) [1], and anecodotally, acquisition scores are similar to a truncated exponential distribution ('80/20' rule). Note that this is a very rough approximation—we do not claim that the distribution of acquisition scores really truncated exponential (see also Figure 7c for a visualization).

At the same time, the choice of a Gumbel distribution for the noise is also one of mathematical convenience, in the spirit of a straightforward baseline. For example, the maximum of sets of many other standard distributions, such as the Gaussian distribution, is not analytically tractable. Taking the highest-scoring points from a distribution perturbed with Gumbel noise is equivalent to sampling from a softmax distribution[2] without replacement.

This follows from the Gumbel-Max trick (Gumbel, 1954; Maddison et al., 2014) and, more specifically, the Gumbel-Top-$K$ trick (Kool et al., 2019). We provide a short proof in appendix B.2. Expanding on Maddison et al. (2014):

---

[1]See also the following Math StackExchange thread.
[2]Also known as Boltzmann/Gibbs distribution.

**Proposition 3.1.** *For scores $s_i$, $i \in \{1, \ldots, n\}$, and $k \leq n$ and $\beta > 0$, if we draw $\epsilon_i \sim Gumbel(0; \beta^{-1})$ independently, then $\arg \operatorname{top}_k \{s_i + \epsilon_i\}_i$ is an (ordered) sample without replacement from the categorical distribution* $\operatorname{Categorical}(\exp(\beta\, s_i) / \sum_j \exp(\beta\, s_j), i \in \{1, \ldots, n\})$.

$\beta \geq 0$ is a 'coldness' parameter, which represents the expected rate at which the scores change as more data is acquired. In the spirit of providing a simple and surprisingly effective baseline without hyperparameters, we fix $\beta := 1$. For $\beta \to \infty$, this distribution will converge towards top-$K$ acquisition. Whereas for $\beta \to 0$, it will converge towards uniform acquisition. We examine ablations of $\beta$ in §6.1.

We apply the perturbation to three quantities in the three sampling schemes: the scores themselves, the log scores, and the rank of the scores. Perturbing the log scores assumes that scores are non-negative and uninformative points should be avoided. Perturbing the ranks can be seen as a robustifying assumption that requires the relative scores to be reliable but allows the absolute scores to be unreliable. We summarise the three versions with their associated sampling distributions are in Table 2.

**Soft-Rank Acquisition.** This first variant only relies on the rank order of the scores and makes no assumptions on whether the acquisition scores are meaningful beyond that. It thus uses the *least* amount of information from the acquisition scores. It only requires the *relative score order* to be useful and ignores the *absolute score values*. If the absolute scores provide useful information, we would expect this method to perform worse than the variants below, which make use of the score values. As we will see, this is indeed sometimes the case .

Ranking the scores $s(i; \mathcal{I}^{\text{train}})$ with descending ranks $\{r_i\}_{i \in \mathcal{I}^{\text{pool}}}$ such that $s(r_i; \mathcal{I}^{\text{train}}) \geq s(r_j; \mathcal{I}^{\text{train}})$ for $r_i \leq r_j$ and smallest rank being 1, we sample index $i$ with probability $\operatorname{p}_{\text{softrank}}(i) \propto r_i^{-\beta}$ with coldness $\beta$. This is invariant to the actual scores. We can draw $\epsilon_i \sim \operatorname{Gumbel}(0; \beta^{-1})$ and create a perturbed 'rank'

$$s^{\text{softrank}}(i; \mathcal{I}^{\text{train}}) := -\log r_i + \epsilon_i. \tag{9}$$

Following Proposition 3.1, taking the top-$K$ points from $s^{\text{softrank}}$ is equivalent to sampling without replacement from the rank distribution $\operatorname{p}_{\text{softrank}}(i)$.

**Softmax Acquisition.** The next simplest variant uses the actual scores instead of the ranks. Again, it perturbs the scores by a Gumbel-distributed random variable $\epsilon_i \sim \operatorname{Gumbel}(0; \beta^{-1})$

$$s^{\text{softmax}}(i; \mathcal{I}^{\text{train}}) := s(i; \mathcal{I}^{\text{train}}) + \epsilon_i. \tag{10}$$

However, this makes no assumptions about the semantics of the absolute values of the scores: the softmax function is invariant to constants shifts. Hence, the sampling distribution will only depend on the *relative scores* and not their absolute value.

**Power Acquisition.** For many scoring functions, the scores are non-negative, and a score close to zero means that the sample is not informative in the sense that we do not expect it will improve the model—we do not want to sample it. This is the case with commonly used score functions such as BALD and entropy. BALD measures the expected information gain. When it is zero for a sample, we do not expect anything to be gained from acquiring a label for that sample. Similarly, entropy is upper-bounding BALD, and the same consideration applies. This assumption also holds for other scoring functions such as the standard deviation and variation ratios; see appendix B.1. To take this into account, the last variant models the future log scores as perturbations of the current log score with Gumbel-distributed noise

$$s^{\text{power}}(i; \mathcal{I}^{\text{train}}) := \log s(i; \mathcal{I}^{\text{train}}) + \epsilon_i. \tag{11}$$

By Proposition 3.1, this is equivalent to sampling from a power distribution

$$\operatorname{p}_{power}(i) \propto \left( \frac{1}{s(i; \mathcal{I}^{\text{train}})} \right)^{-\beta}. \tag{12}$$

This may be seen by noting that $\exp(\beta \log s(i; \mathcal{I}^{\text{train}})) = s(i; \mathcal{I}^{\text{train}})^{\beta}$. Importantly, as scores $\to 0$, the (perturbed) log scores $\to -\infty$ and will have probability mass $\to 0$ assigned. This variant takes the absolute scores into account and avoids data points with score 0.

**In Summary.** Given the above considerations, when using BALD, entropy, and other appropriate scoring functions, power acquisition is the most sensible. Thus, we expect it to work best. Indeed, we find this to be the case in the toy experiment on Repeated-MNIST (Kirsch et al., 2019) depicted in Figure 2. However, even soft-rank acquisition works well in practice, suggesting that the choice of score perturbation is not critical for its effectiveness; see also appendix §D for a more in-depth comparison. In the rest of the main paper, we focus on power acquisition, we include results for all methods in §C.

## 4    Related Work

Researchers in active learning (Atlas et al., 1989; Settles, 2010) have identified the importance of *batch* acquisition as well as the failures of top-$K$ acquisition using straightforward extensions of single-sample methods in a range of settings including support vector machines (Campbell et al., 2000; Schohn & Cohn, 2000; Brinker, 2003; Guo & Schuurmans, 2007), GMMs (Azimi et al., 2012), and neural networks (Sener & Savarese, 2018; Kirsch et al., 2019; Ash et al., 2020; Baykal et al., 2021).

Many of these methods aim to introduce structured diversity to batch acquisition that accounts for the *interaction* of the points acquired in the learning process. In most cases, the computational complexity scales poorly with the acquisition size ($K$) or pool size ($M$), for example because of the estimation of joint mutual information (Kirsch et al., 2019), the $\mathcal{O}(KM)$ complexity of using a k-means++ initialisation scheme (Ash et al., 2020), which approximates k-DPP-based batch active learning (Bıyık et al., 2019), or the $\mathcal{O}(M^2 \log M)$ complexity of methods based on $K$-centre coresets (Sener & Savarese, 2018) (although heuristics and continuous relaxations can improve this somewhat). In contrast, we examine simple and efficient stochastic strategies for adapting well-known single-sample acquisition functions to the batch setting. The proposed stochastic strategies are based on observing that acquisition scores would change as new points are added to the acquisition batch and modelling this difference for additional batch samples in the most naive way, using Gumbel noise. The presented stochastic extensions have the same complexity $\mathcal{O}(M \log K)$ as naive top-$K$ batch acquisition, yet outperform it, and they can perform on par with above more complex methods.

For multi-armed bandits, it has been shown that adding noise to the scores, specifically via Thompson sampling, is effective for choosing informative batches (Kalkanli & Özgür, 2021). Similarly, in reinforcement learning, stochastic prioritisation has been employed as *prioritized replay* (Schaul et al., 2016) which may be effective for reasons analogous to those motivating the approach examined in this work.

While stochastic sampling has not been extensively explored for acquisition in deep active learning, most recently it has been used as an auxiliary step in diversity-based active learning methods that rely on clustering as main mechanism (Ash et al., 2020; Citovsky et al., 2021). Kirsch et al. (2019) empirically find that additional noise in the acquisition scores seems to benefit batch acquisition but do not investigate further. Fredlund et al. (2010) suggest modeling single-point acquisition as sampling from a "*query density*" modulated by the (unknown) sample density p($x$) and analyze a binary classification toy problem. Farquhar et al. (2021) propose stochastic acquisition as part of de-biasing actively learned estimators.

Most relevant to this work, and building on Fredlund et al. (2010) and Farquhar et al. (2021), Zhan et al. (2022) propose a stochastic acquisition scheme that is asymptotically optimal. They normalize the acquisition scores via the softmax function to obtain a query density function for unlabeled samples and draw an acquisition batch from it, similar to SoftmaxEntropy. Their method aims to achieve asymptotic optimality for active learning processes by mitigating the impact of bias. In contrast, in this work, we examine multiple stochastic acquisition strategies based on score-based or rank-based distributions and apply these strategies to several single-sample acquisition functions, such as BALD and entropy (and standard deviation, variation ratios, see Figure 11); and we focus on active learning in a (Bayesian) deep learning setting. As such our empirical results and additional proposed strategies can be seen as complementary to their work.

Thus, while stochastic sampling is generally well-known within acquisition functions, to our knowledge, this work is the first[3] to investigate simple stochastic sampling methods entirely as alternatives to naive top-$K$

---

[3] A workshop version was presented at ICML 2021, and the first submission of this work was concurrent to Zhan et al. (2022).

acquisition in (Bayesian) deep active learning and to compare them to more complex approaches in various settings.

## 5    Experiments

In this section, we empirically verify that the presented stochastic acquisition methods (a) outperform top-$K$ acquisition and (b) are competitive with specially designed batch acquisition schemes like BADGE (Ash et al., 2020) and BatchBALD (Kirsch et al., 2019); and are vastly cheaper than these more complicated methods.

To demonstrate the seriousness of the possible weakness of recent batch acquisition methods, we use a range of datasets. These experiments show that the performance of the stochastic extensions is not dependent on the specific characteristics of any particular dataset. Our experiments include computer vision, natural language processing (NLP), and causal inference (in §6.1). We show that stochastic acquisition helps avoid selecting redundant samples on Repeated-MNIST (Kirsch et al., 2019), examine performance in active learning for computer vision on EMNIST (Cohen et al., 2017), MIO-TCD (Luo et al., 2018), Synbols (Lacoste et al., 2020), and CLINC-150 (Larson et al., 2019) for intent classification in NLP. MIO-TCD is especially close to real-world datasets in size and quality. In appendix C.5, we further investigate edges cases using the Synbols dataset under different types of biases and noise.

Here, we consider both BALD and predictive entropy as scoring functions. We examine other scoring functions on Repeated-MNIST in appendix C.2.1 and observe similar results. For the sake of legible figures, we focus on power acquisition in this section, as it fits BALD and entropy best: the scores are non-negative, and zero scores imply uninformative samples. We show that all three methods (power, softmax, softrank) perform similarly in appendix D.

We are not always able to compare to BADGE and BatchBALD because of computational limitations of those methods. BatchBALD is computationally infeasible for large acquisition sizes ($> 10$) because of time constraints, cf. Table 1. When possible, we use BatchBALD with acquisition size 5 as baseline. Similarly, BADGE runs out of memory for large dataset sizes, such as EMNIST 'ByMerge' with 814,255 examples.

Figures interpolate linearly between available points, and we show 95% confidence intervals.

**Experimental Setup & Compute.** We document the experimental setup and model architectures in detail in appendix C.1. Our experiments used about 25,000 compute hours on Titan RTX GPUs.

**Runtime Measurements.** We emphasize that the stochastic acquisition strategies are much more computationally efficient compared to specialised batch-acquisition approaches like BADGE and BatchBALD. Runtimes, shown in Table 1, are essentially identical for top-$K$ and the stochastic versions. Both are orders of magnitude faster than BADGE and BatchBALD even for small batches. Unlike those methods, stochastic acquisition scales *linearly* in pool size and *logarithmically* in acquisition size. Runtime numbers do not include the cost of retraining models (identical in each case). The runtimes for top-$K$ and stochastic acquisition appear constant over $K$ because the execution time is dominated by fixed-cost memory operations. The synthetic dataset used for benchmarking has 4,096 features, 10 classes, and 10,000 pool points.

**Repeated-MNIST.** Repeated-MNIST (Kirsch et al., 2019) duplicates MNIST a specified number of times and adds Gaussian noise to prevent perfect duplicates. Redundant data are incredibly common in industrial applications but are usually removed from standard benchmark datasets. The controlled redundancies in the dataset allow us to showcase pathologies in batch acquisition methods. We use an acquisition size of 10 and 4 dataset repetitions.

Figure 2 shows that PowerBALD outperforms top-$K$ BALD. While much cheaper computationally, cf. Table 1, PowerBALD also outperforms BADGE and even performs on par with BatchBALD. For BatchBALD, we use an acquisition size of 5, and for BADGE of 20. Note that BatchBALD performs better for smaller acquisition sizes while BADGE (counterintuitively) performs better for larger ones; see Figure 10 in the appendix for an ablation. BatchBALD, BALD, and the stochastic variants all become equivalent for acquisition size 1 when individual points are sampled, which performs best Kirsch et al. (2019).
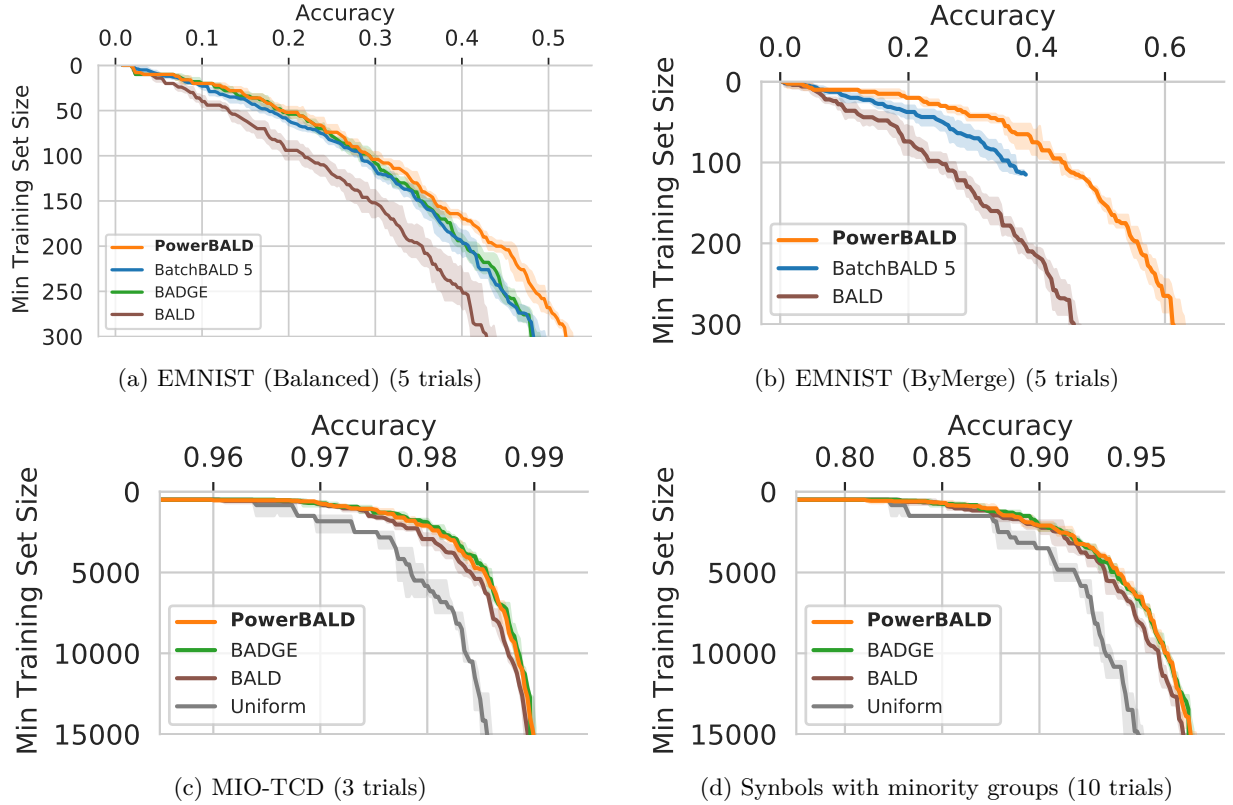
Figure 4: *Performance on various datasets.* BatchBALD took infeasibly long on these datasets & acquisition sizes. **(a)** *EMNIST 'Balanced':* On 132k samples, PowerBALD (acq. size 10) outperforms BatchBALD (acq. size 5) and BADGE (acq. size 40). **(b)** *EMNIST 'ByMerge':* On 814k samples, PowerBALD (acq. size 10) outperforms BatchBALD (acq. size 5). BADGE (not shown) OOM'ed, and BatchBALD took $> 12$ days for 115 acquisitions. **(c)** *MIO-TCD:* PowerBALD performs better than BALD and on par with BADGE (all acq. size 100). **(d)** *Synbols with minority groups:* PowerBALD performs on par with BADGE (all acq. size 100).

**Computer Vision: EMNIST.** EMNIST (Cohen et al., 2017) contains handwritten digits and letters and comes with several splits: we examine the 'Balanced' split with 131,600 samples in Figure 4a[4] and the 'ByMerge' split with 814,255 samples in Figure 4b. Both have 47 classes. We use an acquisition size of 5 for BatchBALD, of 40 for BADGE, and of 10 otherwise.

We see that the stochastic methods outperform BatchBALD on it and both BADGE and BatchBALD on 'Balanced' (Figure 4a). They do not have any issues with the huge pool set in 'ByMerge' (Figure 4b). In the appendix, Figures 24 and 25 show results for all three stochastic extensions, and Figure 16 shows an ablation of different acquisition batch sizes for BADGE. For 'ByMerge', BADGE ran out of memory on our machines, and BatchBALD took more than 12 days for 115 acquisitions when we halted execution.

**Computer Vision: MIO-TCD.** The Miovision Traffic Camera Dataset (MIO-TCD) (Luo et al., 2018) is a vehicle classification and localisation dataset with 648,959 images designed to exhibit realistic data characteristics like class imbalance, duplicate data, compression artefacts, varying resolution (between 100 and 2,000 pixels), and uninformative examples; see Figure 9 in the appendix. As depicted in Figure 4c, PowerBALD performs better than BALD and essentially matches BADGE despite being much cheaper to compute. We use an acquisition size of 100 for all methods.

**Computer Vision: Synbols.** Synbols (Lacoste et al., 2020) is a character dataset generator which can demonstrate the behaviour of batch active learning under various edge cases (Lacoste et al., 2020; Branchaud-

---

[4]This result exactly reproduces BatchBALD's trajectory in Figure 7 from Kirsch et al. (2019).

Charron et al., 2021). In Figure 4d, we evaluate PowerBALD on a dataset with minority character types and colours. PowerBALD outperforms BALD and matches BADGE. Further details as well as an examination of the 'spurious correlation' and 'missing synbols' edge cases (Lacoste et al., 2020; Branchaud-Charron et al., 2021) can be found in appendix C.5.

**Natural Language Processing: CLINC-150.** We perform intent classification on CLINC-150 (Larson et al., 2019), which contains 150 intent classes plus an out-of-scope class. This setting captures data seen in production for chatbots. We fine-tune a pretrained DistilBERT model from HuggingFace (Wolf et al., 2020) on CLINC-150 for 5 epochs with Adam as optimiser. In appendix C.6, we see that PowerEntropy shows strong performance compared to Entropy and performs almost on par with BADGE. This demonstrates that our technique is domain independent and can be easily reused for other tasks.

**In Summary.** We have verified that stochastic acquisition functions outperform top-$K$ batch acquisition in several different settings and perform on par with more complex methods such as BADGE or BatchBALD. Moreover, we refer the reader to Jesson et al. (2021), Murray et al. (2021), Tigas et al. (2022), Holmes et al. (2022) for additional works that use the proposed stochastic acquisition functions in this paper and provide further empirical validation.

## 6 Further Investigations

In this section, we examine and validate assumptions about the underlying score dynamics by examining the scores across acquisitions. We further hypothesise about when top-$K$ acquisition is the most detrimental to active learning.

**Acquisition Asymptotics of Bayesian Models.** For well-specified and well-defined Bayesian parametric models, the posterior distribution of the model parameters converges to the true parameters as the number of data points increases (Van der Vaart, 2000).

For such models and assuming that the predictions are independent given the model parameters, the total correlation between the predictions decreases as the number of training points increases, as the posterior distribution of the model parameters becomes more concentrated around the true parameters:

$$\text{TC}[Y_1, \ldots, Y_K \mid x_1, \ldots, x_K, \mathcal{D}^{\text{train}}] \to 0 \quad \text{as} \quad |\mathcal{D}^{\text{train}}| \to \infty. \tag{13}$$

This can be proved by noting that in the finite data limit, the posterior parameter distribution converges to the true model parameters, and the marginal distribution then factorizes. This means that the predictions become more independent as the number of training points increases and fully independent in the infinite data limit.

The total correlation is defined as:

$$\text{TC}[Y_1, \ldots, Y_K \mid x_1, \ldots, x_K, \mathcal{D}^{\text{train}}] := \underbrace{\sum_i \text{H}[Y_i \mid x_i, \mathcal{D}^{\text{train}}]}_{\text{top-}K \text{ Entropy}} - \underbrace{\text{H}[Y_1, \ldots, Y_K \mid x_1, \ldots, x_K, \mathcal{D}^{\text{train}}].}_{\text{'Batch Entropy'}}, \tag{14}$$

We can also write the total correlation as difference between top-$K$ BALD and BatchBALD:

$$\text{TC}[Y_1, \ldots, Y_K \mid x_1, \ldots, x_K, \mathcal{D}^{\text{train}}] = \underbrace{\sum_i \text{I}[Y_i; \Omega \mid x_i, \mathcal{D}^{\text{train}}]}_{\text{top-}K \text{ BALD}} - \underbrace{\text{I}[Y_1, \ldots, Y_K; \Omega \mid x_1, \ldots, x_K, \mathcal{D}^{\text{train}}]}_{\text{BatchBALD}}. \tag{15}$$

As the total correlation converges to 0, the top-$K$ BALD term (first term) becomes equal to the BatchBALD term (the second term on the right side), and the same happens for top-$K$ entropy and 'BatchEntropy', which we similarly define.

Thus, for well-specified and well-defined Bayesian parametric models, the top-$K$ acquisition functions will eventually become equivalent to the BatchBALD and 'BatchEntropy' acquisition functions as the number of training points increases. This tells us that top-$K$ acquisition is the most detrimental to active learning in the
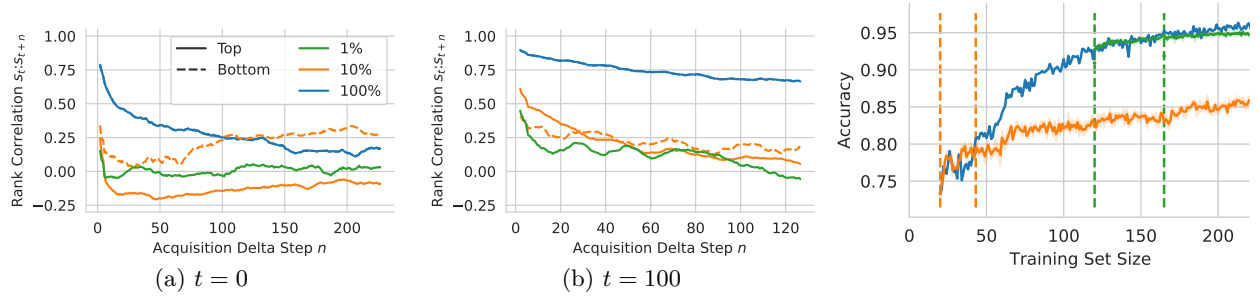
(a) $t = 0$



(b) $t = 100$



Figure 5: *Rank correlations for BALD scores on MNIST between the initial scores and later scores of the top- or bottom-scoring 1%, 10% and 100% of test points (smoothed with a size-10 Parzen window).* Rank-orders decorrelate faster for the most informative samples and in the early stages of training. The top-1% scorers' ranks *anti-correlate* after roughly 40 (100) acquisitions unlike the bottom-1%. Later in training, the acquisition scores stay more strongly correlated. This suggests *the acquisition size could be increased later in training.*

Figure 6: *Top-K acquisition hurts less later in training (BALD on MNIST).* At $t \in \{20, 100\}$ (blue), we keep acquiring samples using the BALD scores from those two steps. At $t = 20$ (orange), the model performs well for $\approx 20$ acquisitions; at $t = 120$ (green), for $\approx 50$; see §6.

earlier stages of learning, when the total correlation between the predictions is still high. This is consistent with our empirical results below ('Increasing Top-$K$ Analysis').

At the same time, as the number of training points increases and the model parameters concentrate, the expected information gain (BALD) also decreases. The mutual information with a deterministic variable is always 0, and thus:

$$ \mathrm{I}[Y; \Omega \mid x, \mathcal{D}^{\mathrm{train}}] \to 0 \quad \text{as} \quad |\mathcal{D}^{\mathrm{train}}| \to \infty. \tag{16} $$

This asymptotic behavior is a trivial but important result, as it tells us that the expected information gain (BALD) will eventually become uninformative as the number of training points increases and no better than random acquisition, and the important question is: when? Given that we only have noisy estimators, this determines until when active learning is of use compared to random acquisition.

Many different active learning methods that are considered non-Bayesian nevertheless approximate the expected information gain or the expected predictive information (Kirsch & Gal, 2022; Smith et al., 2023), which is an expected total correlation. Hence, the considerations apply to those methods, too.

Finally, we can also note that estimators like in BatchBALD which use Monte-Carlo samples of the parameters are bounded by the log of the number of Monte-Carlo samples, which means that they can become uninformative rather quickly: For example, BatchBALD greedily chooses the $t$-th acquisition samples in the batch by maximizing $\mathrm{I}[Y; \Omega \mid x, Y_{t-1}, x_{t-1}, \ldots, Y_1, x_1 \mathcal{D}^{\mathrm{train}}]$ for remaining pool samples $x$. As we have:

$$ \log M \geq \mathrm{I}[Y_1, \ldots, Y_K; \Omega \mid x_1, \ldots, x_K, \mathcal{D}^{\mathrm{train}}] = \sum_{i=1}^{K} \mathrm{I}[Y; \Omega \mid x, Y_y, x_y, \ldots, Y_1, x_1 \mathcal{D}^{\mathrm{train}}], \tag{17} $$

the estimator for $\mathrm{I}[Y_K; \Omega \mid x_K, Y_{K-1}, x_{K-1}, \ldots, Y_1, x_1 \mathcal{D}^{\mathrm{train}}] \to 0$ for increasing $K$ as the overall estimate is bounded by $\log M$. For $M = 100$ parameter samples, we have $\log_{10} M = 2$, and thus BatchBALD could become uninformative after only 2 acquisitions for a classification problem with 10 classes—if there were at least two maximally diverse and uncorrelated points in the pool set.

**Rank Correlations Across Acquisitions.** In Section 3, we made the following assumptions: (1) the acquisition scores $s_t$ at step $t$ are a proxy for scores $s_{t'}$ at step $t' > t$; (2) the larger $t' - t$ is, the worse a proxy $s_t$ is for $s'_t$; (3) this effect is the largest for the most informative points.

We demonstrate these empirically by examining the Spearman rank correlation between scores during acquisition. Specifically, we train a model for $n$ steps using BALD as single-point acquisition function. We compare the rank order at each step to the starting rank order at step $t$. To denoise the rankings across $n$,

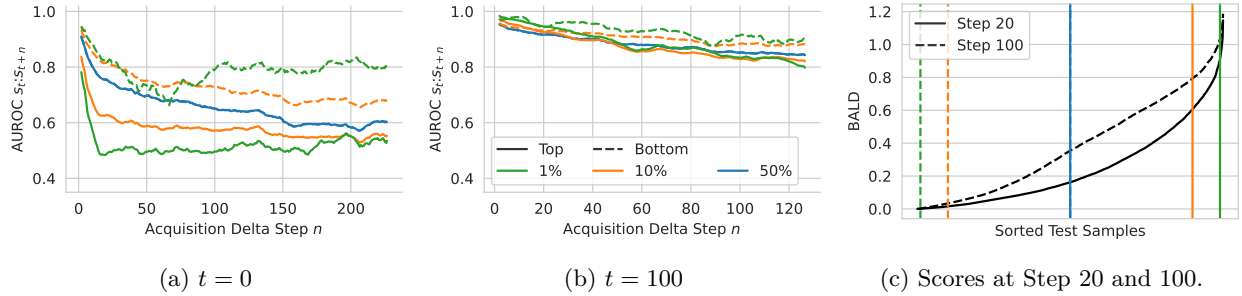(a) $t = 0$            (b) $t = 100$           (c) Scores at Step 20 and 100.

Figure 7: *AUROCs for BALD scores on MNIST between the initial scores and later scores of the top- or bottom-scoring 1%, 10% and 50% of test points (smoothed with a size-10 Parzen window).* AUROC between original points as 'ground truth' and later scores as predictors. This is equivalent to the probablity that the acquisition score at $n$ for a point in $t = 0$'s top or bottom 1%, etc. is larger than points outside. This tells us how likely other points outside the batch have higher acquisition scores. This ignores the ranking of points otherwise. **(a, b)** Points in the top quantiles are superseded by other points in the top quantiles in the later acquisitions to a large degree. This is much more pronounced early in the training than later. The bottom quantiles are more stable. **(c)** The overall score distributions at steps $t = 0, 100$ are visualized and the relevant top and bottom quantiles are marked.

we smooth the rank correlations with a Parzen window of size 10 and to reduce the effect of noise to the rank order, we round all scores to 2 decimal places. This especially removes unimportant rank changes for points with low scores around 0.

Figure 1 shows that acquisition scores become less correlated as more points are acquired. Figure 5a shows this in more detail for the top and bottom 1%, 10% or 100% of scorers of the test set across acquisitions starting at step $t = 0$ for a model initialised with 20 points. The top-10% scoring points (solid green) quickly become uncorrelated across acquisitions and even become *anti-correlated*. In contrast, the points overall (solid blue) correlate well over time (although they have a much weaker training signal on average). This result supports all three of our hypotheses.

At the same time, we see that as training progresses and we converge towards the best model, the order of scores becomes more stable across acquisitions. In Figure 5b the model begins with 120 points ($t = 100$), rather than 20 ($t = 0$). Here, the most informative points are less likely to change their rank—even the top-1% ranks do not become *anti-correlated*, only de-correlated. Thus, we hypothesise that further in training, we might be able to choose larger $K$.

To provide a different analysis, we also consider the more direct question in Figure 7 of how likely other samples have higher acquisition scores at $t + n$ than the top samples from $t$ for different quantiles (1%, 10%, 50%) of the test set. As a sanity check, we also examine the bottom quantiles. This is equivalent to computing the *AUROC* between the original points as 'ground truth' and later scores as predictors. Specifically, we set up a binary classification with the top or bottom 1%, 10% or 50% of the test set as positive and the rest as negative. These results match the previous ones and provide another validation for the mentioned assumptions.

**Increasing Top-$K$ Analysis.** Another way to investigate the effect of top-$K$ selection is to freeze the acquisition scores during training and then continue single-point 'active learning' as if those were the correct scores. Comparing this to the performance of regular active learning with updated single-point scores allows us to examine how well earlier scores perform as proxies for later scores. We perform this toy experiment on MNIST, showing that freezing scores early on greatly harms performance while doing it later has only a small effect (Figure 6). For frozen scores at a training set size of 20 (73% accuracy, $t = 0$), the accuracy matches single-acquisition BALD up to a training set size of roughly 40 (dashed orange lines) before diverging to a lower level. But when freezing the scores of a more accurate model, at a training set size of 120 labels (93% accuracy, $t = 100$), selecting the next fifty points according to those frozen scores performs indistinguishably
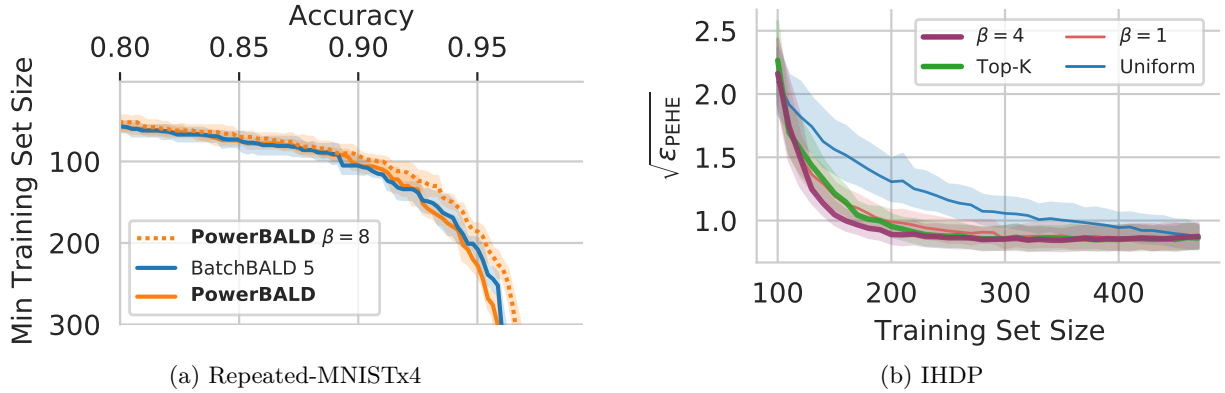
(a) Repeated-MNISTx4

(b) IHDP

Figure 8: *Effect of changing $\beta$*. **(a)** *Repeated-MNISTx4 (5 trials):* PowerBALD outperforms BatchBALD for $\beta = 8$. **(b)** *IHDP (400 trials):* At high temperature ($\beta = 0.1$), CausalBALD with power acquisition is like random acquisition. As the temperature decreases, the performance improves (lower $\sqrt{\epsilon_{\mathrm{PEHE}}}$), surpassing top-$K$ acquisition.

from step-by-step acquisition (dashed green lines). This result shows that top-$K$ acquisition hurts less later in training but can negatively affect performance at the beginning of training.

These observations lead us to ask whether we could dynamically change the acquisition size: with smaller acquisition batches at the beginning and larger ones towards the end of active learning. We leave the exploration of this for future work.

## 6.1 Ablation: Changing $\beta$

So far, we have set $\beta = 1$ in the spirit of examining a simple baseline without additional hyperparameters. The results above show that this already works well and matches the performance of much more expensive methods, raising questions about their value. In addition, however, tuning $\beta$ may be able to further improve performance. In the following, we show that other values of $\beta$ can yield even higher performance on Repeated-MNIST and when estimating causal treatment effects; we provide additional results in appendix E.

**Repeated-MNIST.** In Figure 8a, we see that for PowerBALD the best-performing value, $\beta = 8$, outperforms BatchBALD.

**Causal Treatment Effects: Infant Health Development Programme.** Active learning for Conditional Average Treatment Effect (CATE) estimation Heckman et al. (1997; 1998); Hahn (1998); Abrevaya et al. (2015) on data from the Infant Health and Development Program (IHDP) estimates the causal effect of treatments on an infant's health from observational data. Statistical estimands of the CATE are obtainable from observational data under certain assumptions. Jesson et al. (2021) show how to use active learning to acquire data for label-efficient estimation. Among other subtleties, this prioritises the data for which matched treated/untreated pairs are available.

We follow the experiments of Jesson et al. (2021) on both synthetic data and the semi-synthetic IHDP dataset (Hill, 2011), a commonly used benchmark for causal effects estimation. In Figure 8b we show that power acquisition performs significantly better than both top-$K$ and uniform acquisition, using an acquisition size of 10 in all cases with further. We provide additional results on semi-synthetic data in appendix E.2. Note that methods such as BADGE and BatchBALD are not well-defined for causal-effect estimation, while our approach remains applicable and is effective when fine-tuning $\beta$.

Performance on these tasks is measured using the expected *Precision in Estimation of Heterogeneous Effect (PEHE)* (Hill, 2011) such that $\sqrt{\epsilon_{\mathrm{PEHE}}} = \sqrt{\mathbb{E}[(\widetilde{\tau}(\mathbf{X}) - \tau(\mathbf{X}))^2]}$ (Shalit et al., 2017) where $\widetilde{\tau}$ is the estimated CATE and $\tau$ is CATE (i.e. a form of RMSE).

**Limitations.** Although we highlight the possibility for future work to adapt $\beta$ to specific datasets or score functions, our aim is not to offer a practical recipe for this to practitioners. Our focus is on showing how even the simplest form of stochastic acquisition already raises questions for some recent more complex methods.

## 7    Discussion & Conclusion

We have demonstrated a surprisingly effective and efficient baseline for batch acquisition in active learning. Our stochastic method is orders of magnitude faster than sophisticated batch-acquisition strategies like BADGE and BatchBALD while retaining comparable performance in many settings. Compared to the flawed top-$K$ batch acquisition heuristic, it is never worse: we see no reason to continue using top-$K$ acquisition.

Importantly, our work raises serious questions about these current methods. If they fail to outperform such a simple baseline in a wide range of settings, do they model the interaction between points sufficiently well? If so, are the scores themselves unreliable? We call on future work in batch active learning to at least demonstrate that it can outperform our simple strategy.

At the same time, our framework opens doors for improved methods. Although our stochastic model is put forward for its computational and mathematical simplicity, future work could explore more sophisticated modelling of the predicted score changes that take the current model and dataset into account. In its simplest form, this might mean adapting the temperature of the acquisition distribution to the dataset or estimating it online. Our experiments also highlight that the acquisition size could be dynamic, with larger batch sizes acceptable later in training.

## References

Jason Abrevaya, Yu-Chin Hsu, and Robert P Lieli. Estimating conditional average treatment effects. *Journal of Business
Economic Statistics*, 33(4):485–505, 2015.

Jordan T. Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep batch active learning by diverse, uncertain gradient lower bounds. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL https://openreview.net/forum?id=ryghZJBKPS.

Jordan T. Ash, Surbhi Goel, Akshay Krishnamurthy, and Sham M. Kakade. Gone fishing: Neural active learning with fisher embeddings. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, 2021.

Parmida Atighehchian, Frédéric Branchaud-Charron, and Alexandre Lacoste. Bayesian active learning for production, a systematic study and a reusable library. *arXiv preprint arXiv:2006.09916*, 2020.

Les E. Atlas, David A. Cohn, and Richard E. Ladner. Training connectionist networks with queries and selective sampling. In *Advances in Neural Information Processing Systems 2, [NIPS Conference, Denver, Colorado, USA, November 27-30, 1989]*, 1989.

Javad Azimi, Alan Fern, Xiaoli Zhang Fern, Glencora Borradaile, and Brent Heeringa. Batch active learning via coordinated matching. In *Proceedings of the 29th International Conference on Machine Learning, ICML 2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012*. icml.cc / Omnipress, 2012. URL http://icml.cc/2012/papers/607.pdf.

Cenk Baykal, Lucas Liebenwein, Dan Feldman, and Daniela Rus. Low-regret active learning. *arXiv preprint arXiv:2104.02822*, 2021.

Sarah Bird, Miro Dudík, Richard Edgar, Brandon Horn, Roman Lutz, Vanessa Milan, Mehrnoosh Sameki, Hanna Wallach, and Kathleen Walker. Fairlearn: A toolkit for assessing and improving fairness in AI. Technical Report MSR-TR-2020-32, Microsoft, 2020.

Erdem Bıyık, Kenneth Wang, Nima Anari, and Dorsa Sadigh. Batch active learning using determinantal point processes. *arXiv preprint arXiv:1906.07975*, 2019.

Frédéric Branchaud-Charron, Parmida Atighehchian, Pau Rodríguez, Grace Abuhamad, and Alexandre Lacoste. Can active learning preemptively mitigate fairness issues? *ICLR Workshop on Responsable AI*, 2021.

Klaus Brinker. Incorporating diversity in active learning with support vector machines. In Tom Fawcett and Nina Mishra (eds.), *Machine Learning, Proceedings of the Twentieth International Conference (ICML 2003), August 21-24, 2003, Washington, DC, USA*, pp. 59–66. AAAI Press, 2003. URL http://www.aaai.org/Library/ICML/2003/icml03-011.php.

Colin Campbell, Nello Cristianini, and Alexander J. Smola. Query learning with large margin classifiers. In Pat Langley (ed.), *Proceedings of the Seventeenth International Conference on Machine Learning (ICML 2000), Stanford University, Stanford, CA, USA, June 29 - July 2, 2000*, pp. 111–118. Morgan Kaufmann, 2000.

Gui Citovsky, Giulia DeSalvo, Claudio Gentile, Lazaros Karydas, Anand Rajagopalan, Afshin Rostamizadeh, and Sanjiv Kumar. Batch active learning at scale. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, 2021.

Gregory Cohen, Saeed Afshar, Jonathan Tapson, and André van Schaik. EMNIST: an extension of MNIST to handwritten letters. *arXiv preprint arXiv:1702.05373*, 2017.

Sebastian Farquhar, Yarin Gal, and Tom Rainforth. On statistical bias in active learning: How and when to fix it. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*, 2021.

Richard Fredlund, Richard M. Everson, and Jonathan E. Fieldsend. A bayesian framework for active learning. In *International Joint Conference on Neural Networks, IJCNN 2010, Barcelona, Spain, 18-23 July, 2010*, 2010.

Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1183–1192. PMLR, 2017.

Emil Julius Gumbel. *Statistical theory of extreme values and some practical applications: a series of lectures*, volume 33. US Government Printing Office, 1954.

Yuhong Guo and Dale Schuurmans. Discriminative batch mode active learning. In John C. Platt, Daphne Koller, Yoram Singer, and Sam T. Roweis (eds.), *Advances in Neural Information Processing Systems 20, Proceedings of the Twenty-First Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 3-6, 2007*, pp. 593–600. Curran Associates, Inc., 2007. URL https://proceedings.neurips.cc/paper/2007/hash/ccc0aa1b81bf81e16c676ddb977c5881-Abstract.html.

Jinyong Hahn. On the role of the propensity score in efficient semiparametric estimation of average treatment effects. *Econometrica*, pp. 315–331, 1998.

James J Heckman, Hidehiko Ichimura, and Petra E Todd. Matching as an econometric evaluation estimator: Evidence from evaluating a job training programme. *The review of economic studies*, 64(4):605–654, 1997.

James J Heckman, Hidehiko Ichimura, and Petra Todd. Matching as an econometric evaluation estimator. *The review of economic studies*, 65(2):261–294, 1998.

Jennifer L Hill. Bayesian nonparametric modeling for causal inference. *Journal of Computational and Graphical Statistics*, 20(1):217–240, 2011.

Geoff Holmes, Eibe Frank, Dale Fletcher, and Corey Sterling. Efficiently correcting machine learning: considering the role of example ordering in human-in-the-loop training of image classification models. In *IUI 2022: 27th International Conference on Intelligent User Interfaces, Helsinki, Finland, March 22 - 25, 2022*, 2022.

Neil Houlsby, Ferenc Huszar, Zoubin Ghahramani, and Máté Lengyel. Bayesian active learning for classification and preference learning. *arXiv preprint arXiv:1112.5745*, 2011.

Andrew Jesson, Panagiotis Tigas, Joost van Amersfoort, Andreas Kirsch, Uri Shalit, and Yarin Gal. Causal-bald: Deep bayesian active learning of outcomes to infer treatment-effects from observational data. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, 2021.

Cem Kalkanli and Ayfer Özgür. Batched thompson sampling. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, 2021.

Andreas Kirsch and Yarin Gal. Unifying approaches in active learning and active sampling via fisher information and information-theoretic quantities. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL https://openreview.net/forum?id=UVDAKQANOW.

Andreas Kirsch, Joost van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 7024–7035, 2019.

Wouter Kool, Herke van Hoof, and Max Welling. Stochastic beams and where to find them: The gumbel-top-k trick for sampling sequences without replacement. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 3499–3508. PMLR, 2019.

Andreas Krause and Daniel Golovin. Submodular function maximization. *Tractability*, 3:71–104, 2014.

Alex Kulesza and Ben Taskar. k-dpps: Fixed-size determinantal point processes. In *Proceedings of the 28th International Conference on Machine Learning, ICML 2011, Bellevue, Washington, USA, June 28 - July 2, 2011*, 2011.

Alexandre Lacoste, Pau Rodríguez López, Frederic Branchaud-Charron, Parmida Atighehchian, Massimo Caccia, Issam Hadj Laradji, Alexandre Drouin, Matt Craddock, Laurent Charlin, and David Vázquez. Synbols: Probing learning algorithms with synthetic datasets. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/0169cf885f882efd795951253db5cdfb-Abstract.html.

Stefan Larson, Anish Mahendran, Joseph J. Peper, Christopher Clarke, Andrew Lee, Parker Hill, Jonathan K. Kummerfeld, Kevin Leach, Michael A. Laurenzano, Lingjia Tang, and Jason Mars. An evaluation dataset for intent classification and out-of-scope prediction. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, 2019.

Zhiming Luo, Frederic Branchaud-Charron, Carl Lemaire, Janusz Konrad, Shaozi Li, Akshaya Mishra, Andrew Achkar, Justin A. Eichel, and Pierre-Marc Jodoin. MIO-TCD: A new benchmark dataset for vehicle classification and localization. *IEEE Trans. Image Process.*, 2018.

Chris J. Maddison, Daniel Tarlow, and Tom Minka. A* sampling. In Zoubin Ghahramani, Max Welling, Corinna Cortes, Neil D. Lawrence, and Kilian Q. Weinberger (eds.), *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pp. 3086–3094, 2014.

Jishnu Mukhoti, Andreas Kirsch, Joost van Amersfoort, Philip HS Torr, and Yarin Gal. Deep deterministic uncertainty: A simple baseline. *arXiv e-prints*, pp. arXiv–2102, 2021.

Chelsea Murray, James Urquhart Allingham, Javier Antorán, and José Miguel Hernández-Lobato. Depth uncertainty networks for active learning. *arXiv preprint arXiv:2112.06796*, 2021.

Jerzy Neyman. edited and translated by dorota m. dabrowska and terrence p. speed (1990). on the application of probability theory to agricultural experiments. essay on principles. section 9. *Statistical Science*, 5(4): 465–472, 1923.

Donald B Rubin. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of educational Psychology*, 66(5):688, 1974.

Donald B Rubin. Randomization analysis of experimental data: The fisher randomization test comment. *Journal of the American Statistical Association*, 75(371):591–593, 1980.

Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. In Yoshua Bengio and Yann LeCun (eds.), *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016. URL http://arxiv.org/abs/1511.05952.

Greg Schohn and David Cohn. Less is more: Active learning with support vector machines. In Pat Langley (ed.), *Proceedings of the Seventeenth International Conference on Machine Learning (ICML 2000), Stanford University, Stanford, CA, USA, June 29 - July 2, 2000*, pp. 839–846. Morgan Kaufmann, 2000.

Jasjeet S Sekhon. The neyman-rubin model of causal inference and estimation via matching methods. *The Oxford handbook of political methodology*, 2:1–32, 2008.

Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL https://openreview.net/forum?id=H1aIuk-RW.

Burr Settles. Active Learning Literature Survey. *Machine Learning*, 2010.

Uri Shalit, Fredrik D. Johansson, and David A. Sontag. Estimating individual treatment effect: generalization bounds and algorithms. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pp. 3076–3085. PMLR, 2017. URL http://proceedings.mlr.press/v70/shalit17a.html.

Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In Yoshua Bengio and Yann LeCun (eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL http://arxiv.org/abs/1409.1556.

Freddie Bickford Smith, Andreas Kirsch, Sebastian Farquhar, Yarin Gal, Adam Foster, and Tom Rainforth. Prediction-oriented bayesian active learning. *arXiv preprint arXiv:2304.08151*, 2023.

Panagiotis Tigas, Yashas Annadani, Andrew Jesson, Bernhard Schölkopf, Yarin Gal, and Stefan Bauer. Interventions, where and how? experimental design for causal models at scale. *arXiv preprint arXiv:2203.02016*, 2022.

Joost van Amersfoort, Lewis Smith, Andrew Jesson, Oscar Key, and Yarin Gal. Improving deterministic uncertainty estimation in deep learning for classification and regression. *arXiv preprint arXiv:2102.11409*, 2021.

Aad W Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.

Sahil Verma and Julia Rubin. Fairness definitions explained. In *Proceedings of the International Workshop on Software Fairness, FairWare@ICSE 2018, Gothenburg, Sweden, May 29, 2018*, 2018.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, EMNLP 2020 - Demos, Online, November 16-20, 2020*, 2020.

Xueying Zhan, Yaowei Wang, and Antoni B. Chan. Asymptotic optimality for active learning processes. In *Uncertainty in Artificial Intelligence, Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence, UAI 2022, 1-5 August 2022, Eindhoven, The Netherlands*, 2022.

# A  Ethical impact

We do not foresee any ethical risks related to this work. Insofar as our sampling method reduces computational costs, applications might benefit from reduced resource consumption. Our method appears to be better than or as good as alternatives on evaluations examining the ability to learn from data with under-represented groups and on evaluations that measure the difference between performance for the most- and least-represented groups, which may aid algorithmic fairness (see C.5).

# B  Method

## B.1  Other scoring functions

Following Gal et al. (2017), we also examine using variation ratios (least confidence) and standard deviation as scoring functions.

**Variation Ratio.** Also known as *least confidence*, the variation-ratios is the complement of the least-confindent class prediction:

$$s_{\text{variation-ratios}}(i; \mathcal{I}^{\text{train}}) := 1 - \max_y \mathrm{p}(y \mid X = x_i). \tag{18}$$

This scoring function is non-negative and a score of 0 means that the sample is uninformative: a score of 0 means that the respective prediction is one-hot, which means that the expected information gain is also 0 as can be easily verified. Thus, variation ratios matches the intuitions behind power acquisition.

**Standard Deviation.** The standard deviation score function measures the sum of the class probability deviations and is closely related to the BALD scores:

$$s_{\text{std-dev}}(i; \mathcal{I}^{\text{train}}) := \sum_y \sqrt{\mathrm{Var}_{\mathrm{p}(\omega)}[\mathrm{p}(y \mid X = x_i, \omega)]}. \tag{19}$$

This scoring function is also non-negative, and no variance for the predictions implies a zero expected information gain and thus an uninformative sample. Thus, the standard deviation should also perform well with power acquisition.

## B.2  Proof of Proposition 3.1

First, we remind the reader that a random variable $G$ is Gumble distributed $G \sim \text{Gumbel}(\mu; \beta)$ when its cumulative distribution function follows $\mathrm{p}(G \leq g) = \exp(-\exp(-\frac{g-\mu}{\beta}))$.

Furthermore, the Gumbel distribution is closed under translation and positive scaling:

**Lemma B.1.** *Let $G \sim \text{Gumbel}(\mu; \beta)$ be a Gumbel distributed random variable, then:*

$$\alpha G + d \sim \text{Gumbel}(d + \alpha\mu; \alpha\beta). \tag{20}$$

*Proof.* We have $\mathrm{p}(\alpha G + d \leq x) = \mathrm{p}(G \leq \frac{x-d}{\alpha})$. Thus, we have:

$$\mathrm{p}(\alpha G + d \leq x) = \exp(-\exp(-\frac{\frac{x-d}{\alpha} - \mu}{\beta})) \tag{21}$$

$$= \exp(-\exp(-\frac{x - (d + \alpha\mu)}{\alpha\beta})) \tag{22}$$

$$\Leftrightarrow \alpha G + d \sim \text{Gumbel}(d + \alpha\mu; \alpha\beta). \tag{23}$$

$\square$

We can then easily prove Proposition 3.1 using Theorem 1 from Kool et al. (2019), which we present it here slightly reformulated to fit our notation:

(a) A good example in MIOTCD dataset.

(b) An example of duplicated samples in the dataset.

(c) An example of class confusion between motorcycle and bicycle.

(d) An example of heavy compression artefact.

(e) An example of low resolution samples.

Figure 9: *MIO-TCD Dataset* is designed to include common artifacts from production data. The size and quality of the images vary greatly between crops; from high-quality cameras on sunny days to low-quality cameras at night. (a) shows an example of clean samples that can be clearly assigned to a class. (b)(c)(d) and (e) show the different categories of noise. (b) shows an example of many near-duplicates that exist in the dataset. (c) is a good example where the assigned class is subject to interpretation (d) is a sample with heavy compression artefacts and (e) is an example of samples with low resolution which again is considered a hard example to learn for the model.

**Lemma B.2.** *For $k \leq n$, let $I_1^*, \ldots, I_k^* = \arg \text{top}_k \{s_i + \epsilon_i\}_i$ with $\epsilon_i \sim \text{Gumbel}(0;1)$, i.i.d.. Then $I_1^*, \ldots, I_k^*$ is an (ordered) sample without replacement from the $\text{Categorical}\left(\frac{\exp s_i}{\sum_{j \in n} \exp s_j}, i \in \{1, \ldots, n\}\right)$ distribution, e.g. for a realization $i_1^*, \ldots, i_k^*$ it holds that*

$$P\left(I_1^* = i_1^*, \ldots, I_k^* = i_k^*\right) = \prod_{j=1}^k \frac{\exp s_{i_j^*}}{\sum_{\ell \in N_j^*} \exp s_\ell}$$

*where $N_j^* = N \setminus \{i_1^*, \ldots, i_{j-1}^*\}$ is the domain (without replacement) for the j-th sampled element.*

Now, it is easy to prove the proposition:

**Proposition 3.1.** *For scores $s_i$, $i \in \{1, \ldots, n\}$, and $k \leq n$ and $\beta > 0$, if we draw $\epsilon_i \sim \text{Gumbel}(0; \beta^{-1})$ independently, then $\arg \text{top}_k \{s_i + \epsilon_i\}_i$ is an (ordered) sample without replacement from the categorical distribution $\text{Categorical}(\exp(\beta s_i) / \sum_j \exp(\beta s_j), i \in \{1, \ldots, n\})$.*

*Proof.* As $\epsilon_i \sim \text{Gumbel}(0; \beta^{-1})$, define $\epsilon_i' := \beta \epsilon_i \sim \text{Gumbel}(0;1)$. Further, let $s_i' := \beta s_i$. Applying Lemma B.2 on $s_i'$ and $\epsilon_i'$, $\arg \text{top}_k \{s_i' + \epsilon_i'\}_i$ yields (ordered) samples without replacement from the categorical distribution $\text{Categorical}(\frac{\exp(\beta s_i)}{\sum_j \exp(\beta s_j)}, i \in \{1, \ldots, n\})$. However, multiplication by $\beta$ does not change the resulting indices of $\arg \text{top}_k$:

$$\arg \text{top}_k \{s_i' + \epsilon_i'\}_i = \arg \text{top}_k \{s_i + \epsilon_i\}_i, \tag{24}$$

concluding the proof. □

## C   Experiments

### C.1   Experimental setup & compute

Full code for all experiments will be available at `anonymized_github_repo`.

**Frameworks.** We use PyTorch. Repeated-MNIST and EMNIST experiments use PyTorch Ignite. Synbols and MIO-TCD experiments use the BaaL library `https://github.com/baal-org/baal` (Atighehchian et al., 2020). Predictive parity is calculated using FairLearn (Bird et al., 2020). The CausalBALD experiments use `https://github.com/anndvision/causal-bald` (Jesson et al., 2021).

**Compute.** Results shown in Table 1 were run inside Docker containers with 8 CPUs (2.2Ghz) and 32 Gb of RAM. Other experiments were run on similar machines with Titan RTX GPUs. The Repeated-MNIST and EMNIST experiments take about 5000 GPU hours. The MIO, Synbols and CLINC-150 experiments take about 19000 GPU hours. The CausalBALD experiments take about 1000 GPU hours.

**Dataset Licenses.** Repeated-MNIST is based on MNIST which is made available under the terms of the Creative Commons Attribution-Share Alike 3.0 license. The EMNIST dataset is made available as CC0 1.0 Universal Public Domain Dedication. Synbols is a dataset generator. MIO-TCD is made available under the terms of the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. CLINC-150 is made available under the terms of Creative Commons Attribution 3.0 Unported License.

#### C.1.1   Runtime measurements

The synthetic dataset used for benchmarking has 4,096 features, 10 classes, and 10,000 pool points. VGG-16 models (Simonyan & Zisserman, 2015) were used to sample predictions and latent embeddings.

#### C.1.2   Repeated-MNIST

The Repeated-MNIST dataset is also constructed following Kirsch et al. (2019) with duplicated examples from MNIST with isotropic Gaussian noise added to the input images (standard deviation 0.1).

We use the same setup as Kirsch et al. (2019): a LeNet-5-like architecture with ReLU activations instead of tanh and added dropout. The model obtains 99% test accuracy when trained on the full MNIST dataset. Specifically, the model is made up of two blocks of a convolution, dropout, max-pooling, ReLU with 32 and 64 channels and 5x5 kernel size, respectively. As classifier head, a two-layer MLP with 128 hidden units (and 10 output units) is used that includes dropout between the layers. We use a dropout probability of 0.5 everywhere. The model is trained with early stopping using the Adam optimiser and a learning rate of 0.001. We sample predictions using 100 MC-Dropout samples for BALD. Weights are reinitialized after each acquisition step.

#### C.1.3   EMNIST

We follow the setup from (Kirsch et al., 2019) with 20 MC dropout samples. We use a similar model as for Repeated-MNIST but with three blocks instead of two. Specifically, we use 32, 64, and 128 channels and 3x3 kernel size. This is followed by a 2x2 max pooling layer before the classifier head. The classifier head is a two-layer MLP but with 512 hidden units instead of 128. Again, we use dropout probability 0.5 everywhere.

#### C.1.4   Synbols & MIO-TCD

The full list of hyperparameters for the Synbols and MIO-TCD experiments is presented in Table 3. Our experiments are built using the BaaL library (Atighehchian et al., 2020). We compute the predictive parity using FairLearn (Bird et al., 2020). We use VGG-16 model (Simonyan & Zisserman, 2015) trained for 10 epochs using Monte Carlo dropout for acquisition (Gal et al., 2017) with 20 dropout samples.

In Figure 9, we show a set of images with common problems that can be find in MIO-TCD.

Table 3: Hyper-parameters used in Section 5 and C.5

| Hyperparameter | Value |
|---|---|
| Learning rate | 0.001 |
| Optimizer | SGD |
| Weight decay | 0 |
| Momentum | 0.9 |
| Loss function | Crossentropy |
| Training duration | 10 |
| Batch size | 32 |
| Dropout $p$ | 0.5 |
| MC iterations | 20 |
| Query size | 100 |
| Initial set | 500 |

### C.1.5 CLINC-150

We fine-tune a pretrained DistilBERT model from HuggingFace (Wolf et al., 2020) on CLINC-150 for 5 epochs with Adam as optimiser. Estimating epistemic uncertainty in transformer models is an open research question, and hence, we do not report results using BALD and focus on entropy instead.

### C.1.6 CausalBALD

Using the Neyman-Rubin framework (Neyman, 1923; Rubin, 1974; Sekhon, 2008), the CATE is formulated in terms of the potential outcomes, $Y_t$, of treatment levels $t \in \{0, 1\}$. Given observable covariates, $\mathbf{X}$, the CATE is defined as the expected difference between the potential outcomes at the measured value $\mathbf{X} = \mathbf{x}$: $\tau(\mathbf{x}) = \mathbb{E}[Y_1 - Y_0 \mid \mathbf{X} = \mathbf{x}]$. This causal quantity is fundamentally unidentifiable from observational data without further assumptions because it is not possible to observe both $Y_1$ and $Y_0$ for a given unit. However, under the assumptions of consistency, non-interference, ignoreability, and positivity, the CATE is identifiable as the statistical quantity $\widetilde{\tau}(\mathbf{x}) = \mathbb{E}[Y \mid T = 1, \mathbf{X} = \mathbf{x}] - \mathbb{E}[Y \mid T = 0, \mathbf{X} = \mathbf{x}]$ (Rubin, 1980).

Jesson et al. (2021) define BALD acquisition functions for active learning CATE functions from observational data when the cost of acquiring an outcome, y, for a given covariate and treatment pair, $(\mathbf{x}, t)$, is high. Because we do not have labels for $Y_1$ and $Y_0$ for each $(\mathbf{x}, t)$ pair in the dataset, their acquisition function focusses on acquiring data points $(\mathbf{x}, t)$ for which it is likely that a matched pair $(\mathbf{x}, 1 - t)$ exists in the pool data or has already been acquired at a previous step. We follow their experiments on their synthetic dataset with limited positivity and the semi-synthetic IHDP dataset (Hill, 2011). Details of the experimental setup are given in (Jesson et al., 2021), we use their provided code, and implement the power acquisition function.

The settings for causal inference experiments are identical to those used in Jesson et al. (2021), using the IHDP dataset (Hill, 2011). Like them, we use a Deterministic Uncertainty Estimation Model (van Amersfoort et al., 2021), which is initialised with 100 datapoints and acquire 10 datapoints per acquisition batch for 38 steps. The dataset has 471 pool points and a 201 point validation set.

## C.2 Repeated-MNIST



Figure 10: *Repeated-MNIST x4 (5 trials): acquisition size ablation for BADGE.* Acquisition size 20 performs best out of $\{10, 20, 40\}$. Hence, we use that for Figure 2.

**BADGE Ablation.** In Figure 10, we see that BADGE performs best with acquisition size 20 on Repeated-MNISTx4 overall. BADGE 40 and BADGE 20 have the highest final accuracy, cf. BADGE 10 while BADGE 20 performs better than BADGE 40 for small training set sizes.
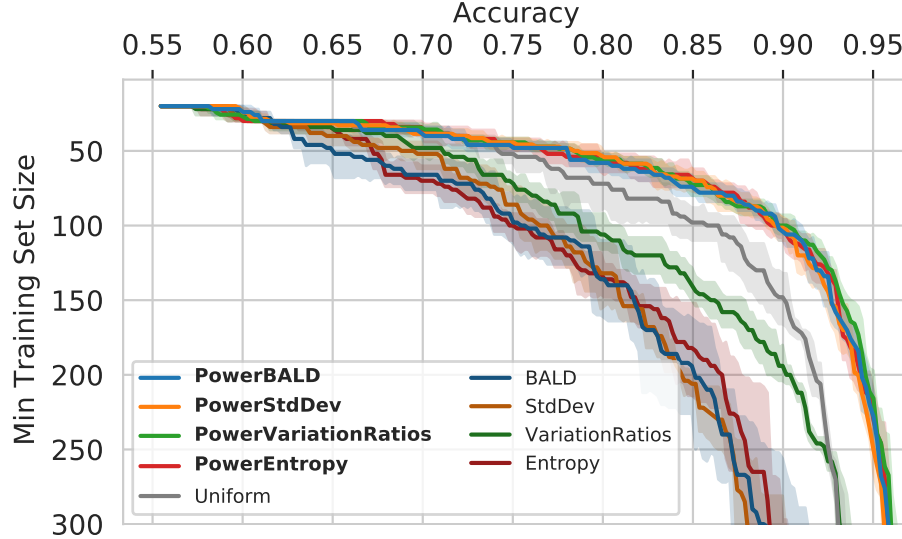
### C.2.1 Other scoring functions



Figure 11: *Repeated-MNIST x4 (5 trials): Performance for other scoring functions.* Entropy, std dev, variation ratios behave like BALD when applying our stochastic sampling scheme.

In Figure 11 shows the performance of other scoring functions than BALD on RepeatedMNIST x4.

### C.2.2   Redundancy ablation



Figure 12: *Repeated-MNIST (5 trials): Performance ablation for different repetition counts.*

In Figure 12, we see the same behaviour in an ablation for different repetition sizes of Repeated-MNIST.

### C.3   MIO-TCD



(a) BALD

(b) Entropy

Figure 13: *MIO-TCD (5 trials).*

In Figure 13, we see that power acquisition performs on par with BADGE with both BALD and entropy as underlying score functions.

### C.4   EMNIST



Figure 14: *EMNIST (Balanced) (5 trials): Performance with BALD.*



Figure 15: *EMNIST (ByMerge) (5 trials): Performance with BALD.*

In Figure 14 and 15, we see that PowerBALD outperforms BALD, BatchBALD, and BADGE.



Figure 16: *EMNIST (Balanced) (5 trials): acquisition size ablation for BADGE.*

**BADGE Ablation.** In Figure 16, we see that BADGE performs similarly with all three acquisition sizes. Acquisition size 10 is the smoothest.

## C.5 Edge cases in Synbols

We use Synbols (Lacoste et al., 2020) to demonstrate the behaviour of batch active learning in artificially constructed edge cases. Synbols is a character dataset generator for classification where a user can specify the type and proportion of bias and insert artefacts, backgrounds, masking shapes, and so on. We selected three datasets with strong biases supplied by Lacoste et al. (2020); Branchaud-Charron et al. (2021) to evaluate our method. The experimental settings are described in appendix C.1.

For these tasks, performance evaluation includes 'predictive parity', also known as 'accuracy difference', which is the maximum difference in accuracy between subgroups—which are, in this case, different coloured characters. This measure is used most widely in domain adaptation and ethics (Verma & Rubin, 2018). We want to maximise the accuracy while minimising the predictive parity.



(a) Accuracy

(b) Predictive parity (**Down and left is better.**)

Figure 17: *Performance on Synbols Spurious Correlations (3 trials) with BALD.* Stochastic acquisition matches BADGE and BALD's predictive parity and performance, which is reassuring as stochastic acquisition functions might be affected by spurious correlations.

**Spurious Correlations.** This dataset includes spurious correlations between character colour and class. As shown in Branchaud-Charron et al. (2021), active learning is especially strong here as characters that do not follow the correlation will be informative and thus selected.

We compare the predictive parity between methods in Fig. 17b. We do not see any significant difference between our method and BADGE or BALD. This is encouraging, as stochastic approaches might select more examples following the spurious correlation and thus have higher predictive parity, but this is not the case.



(a) Accuracy

(b) Predictive parity

Figure 18: *Synbols Minority Groups (3 trials): Performance on BALD.* PowerBALD outperforms BALD and matches BADGE for both accuracy and predictive parity.

**Minority Groups.** This dataset includes a subgroup of the data that is under-represented; specifically, most characters are red while few are blue. As Branchaud-Charron et al. (2021) shows, active learning can improve the accuracy for these groups.

Our stochastic approach lets batch acquisition better capture under-represented subgroups. In Figure 18a, PowerBALD has an accuracy almost identical to that of BADGE, despite being much cheaper, and outperforms BALD. At the same time, we see in Figure 18b that PowerBALD has a lower predictive parity than BALD, demonstrating a fairer predictive distribution given the unbalanced dataset.



Figure 19: BALD

Figure 20: Entropy

Figure 21: *Performance on Synbols Missing Characters (3 trials).* In this dataset with high aleatoric uncertainty, PowerBALD matches BADGE and BALD performance. PowerEntropy significantly outperforms Entropy which confounds aleatoric and epistemic uncertainty.

**Missing Synbols.** This dataset has high aleatoric uncertainty (input noise). Some images are missing information required to make high-probability predictions—these images have shapes randomly occluding the character—so even a perfect model would remain uncertain. Lacoste et al. (2020) demonstrated that entropy is ineffective on this data as it cannot distinguish between aleatoric and epistemic uncertainty (input noise and model uncertainty), while BALD can do so. As a consequence, entropy will unfortunately prefer samples with occluded characters, resulting in degraded active learning performance. For predictive entropy, stochastic acquisition largely corrects the failure of entropy acquisition to account for missing data (Figure 21) although PowerEntropy still underperforms BADGE here. For BALD, we show in Figure 19 in the appendix that, as before, our stochastic method performs on par with BADGE and marginally better than BALD.

## C.6 CLINC-150



Figure 22: *Performance on CLINC-150 (10 trials).* PowerEntropy performs much better than entropy, which only performs marginally better than uniform, and almost on par with BADGE.

In Figure 22, we see that PowerEntropy performs much better than entropy which only performs marginally better than the uniform baseline. PowerEntropy also performs better than BADGE at low training set sizes, but BADGE performs better in the second half. Between $\approx 2300$ and $4000$ samples, BADGE and PowerEntropy perform the same.

# D Comparing Power, Softmax and Soft-Rank

## D.1 Empirical Evidence



Figure 23: *Repeated-MNIST (5 trials): Performance with all three stochastic strategies.*

**Repeated-MNIST.** In Figure 23, power acquisition performs best overall, followed by soft-rank and then softmax.
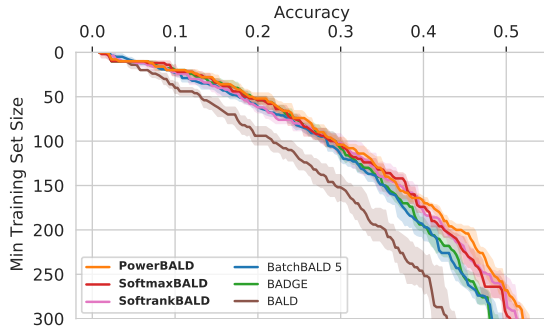


Figure 24: *EMNIST (Balanced) (5 trials): Performance with all three stochastic strategies with BALD.* PowerBALD performs best.
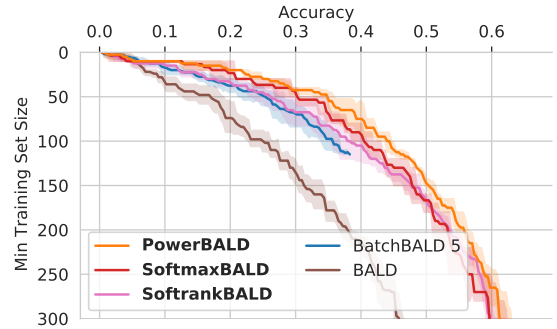


Figure 25: *EMNIST (ByMerge) (5 trials): Performance with all three stochastic strategies with BALD.* PowerBALD performs best.

**EMNIST.** In Figure 24 and 25, we see that PowerBALD performs best, but Softmax- and SoftrankBALD also outperform other methods. BADGE did not run on EMNIST (ByMerge) due to out-of-memory issues and BatchBALD took very long as EMNIST (ByMerge) has more than 800,000 samples.

28

(a) BALD

(b) Entropy

Figure 26: *MIO-TCD (3 trials): Performance with all three stochastic strategies.*

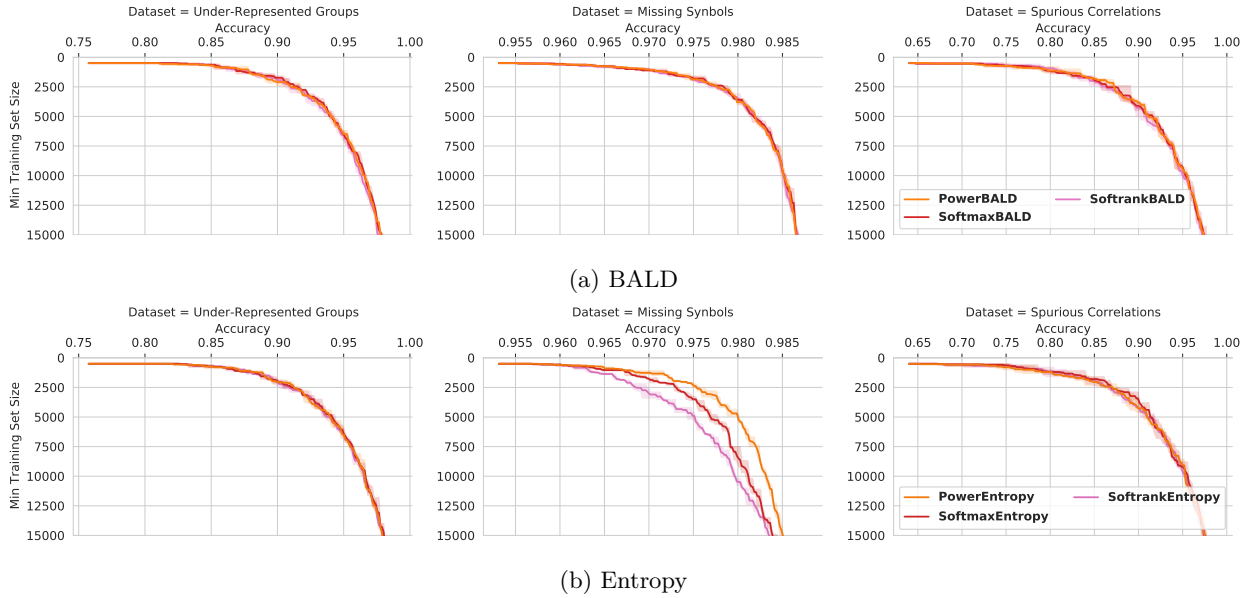**MIO-TCD.** In Figure 26, we see that all three stochastic acquisition methods perform about equally well.



(a) BALD



(b) Entropy

Figure 27: *Synbols edge cases (3 trials): Performance with all three stochastic strategies.*

**Synbols.** In Figure 27, power acquisition seems to perform better overall—mainly due to the performance in Synbols Missing Characters.
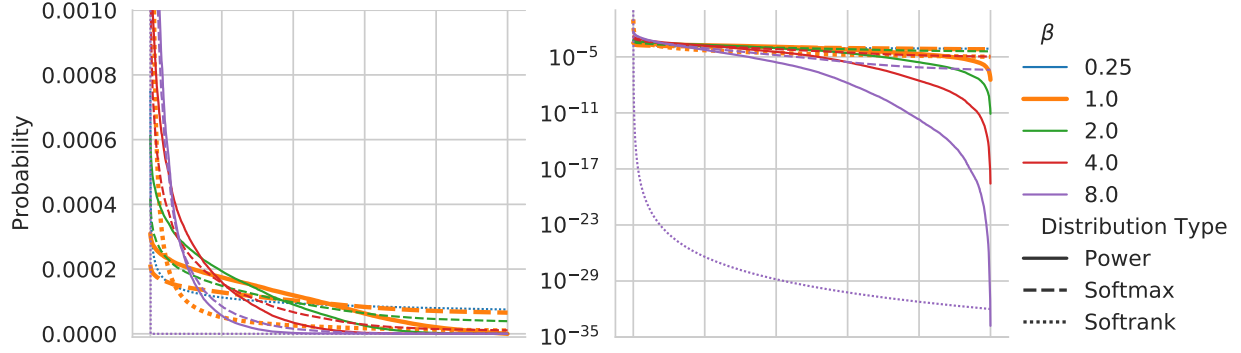
Figure 29: *Score distribution for power and softmax acquisition of BALD scores on MNIST for varying Coldness $\beta$ at $t = 0$.* Linear and log plot over samples sorted by their BALD score. At $\beta = 8$ both softmax and power acquisition have essentially the same distribution for high scoring points (closely followed by the power distribution for $\beta = 4$). This might explain why the coldness ablation shows that these $\beta$ to have very similar AL trajectories on MNIST. Yet, while softmax and power acquisition seem transfer to RMNIST, this is not the case for softrank which is much more sensitive to $\beta$. At the same time, power acquisition avoids low-scoring points more than softmax acquisition.



Figure 28: *CLINC-150 (10 trials): Performance with all three stochastic strategies.*

**CLINC-150.** In Figure 28, all three stochastic methods perform similarly.

### D.2 Investigation

To further examine the three stochastic acquisition variants, we plot their score distributions, extracted from the same MNIST toy example, in Figure 29. Power and softmax acquisition distributions are similar for $\beta = 8$ (power, softmax) and $\beta = 4$ (softmax). This might explain why active learning with these $\beta$ shows similar accuracy trajectories.

We find that power and softmax acquisition are quite insensitive to $\beta$ and thus selecting $\beta = 1$ might generally work quite well.

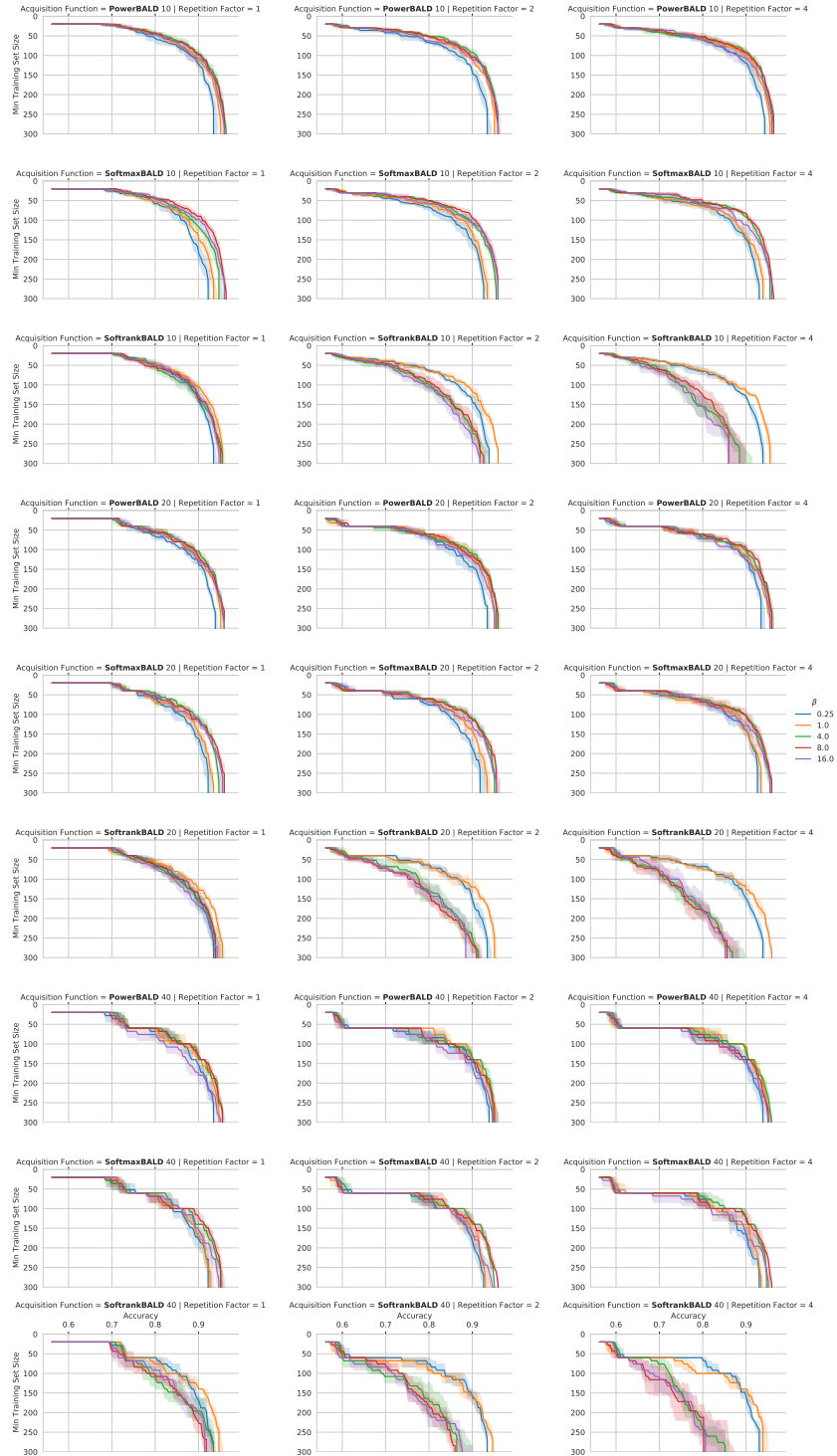# E  Effect of changing $\beta$

## E.1  Repeated-MNIST



Figure 30: *Repeated-MNIST: $\beta$ ablation for *BALD.*
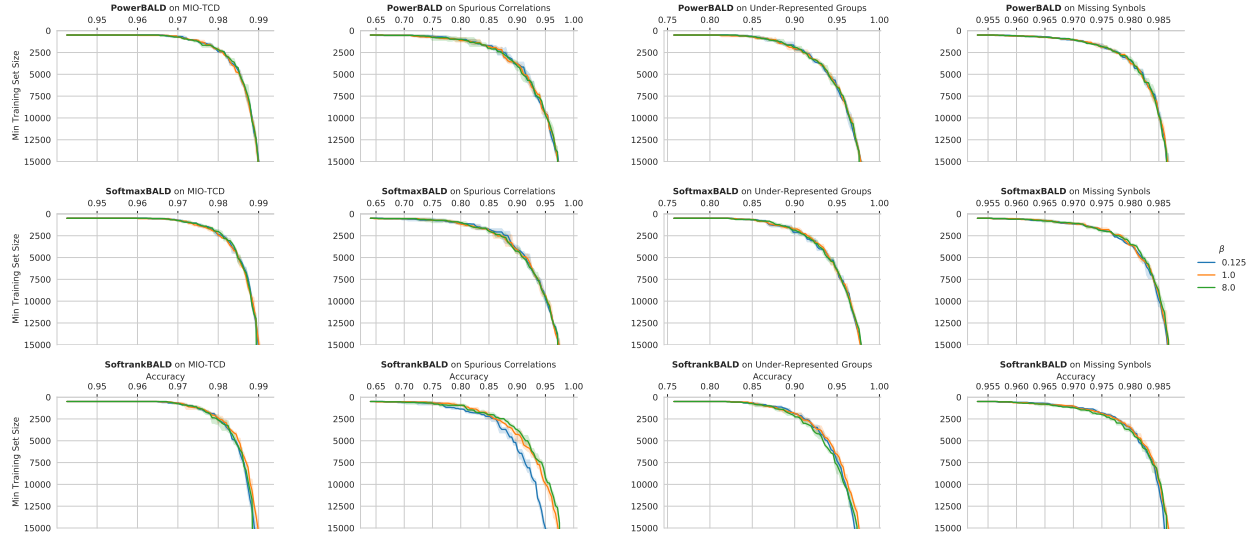
### E.1.1 MIO-TCD and Synbols



Figure 31: *MIO-TCD and Synbols: β ablation for \*BALD.*
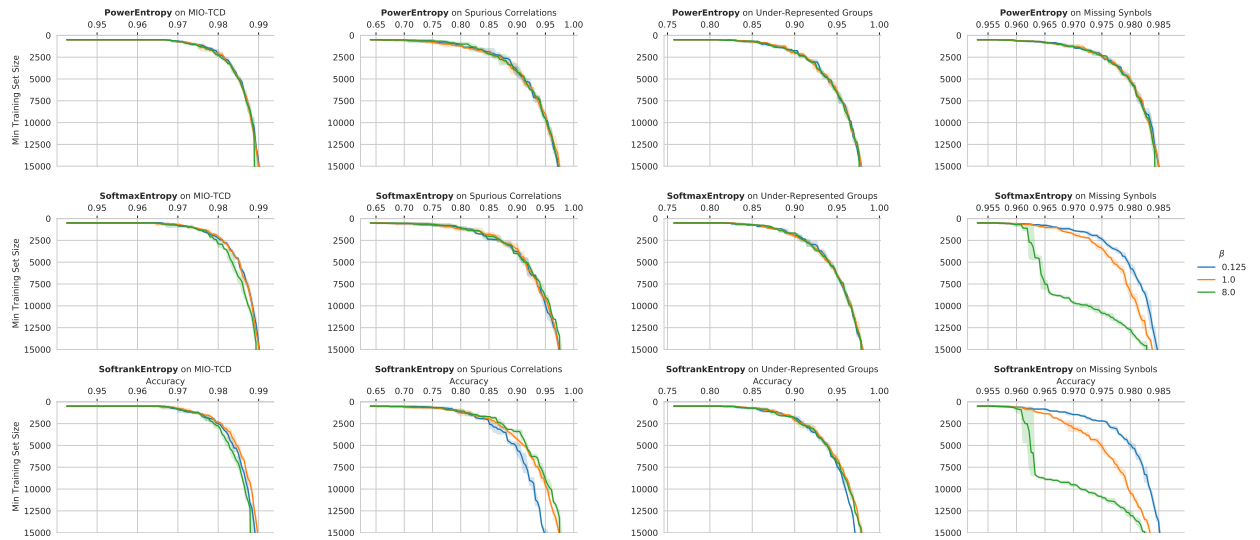


Figure 32: *MIO-TCD and Synbols: β ablation for \*Entropy.*

## E.2 CausalBALD: synthetic dataset



(a) Overall Ablation (Subset)

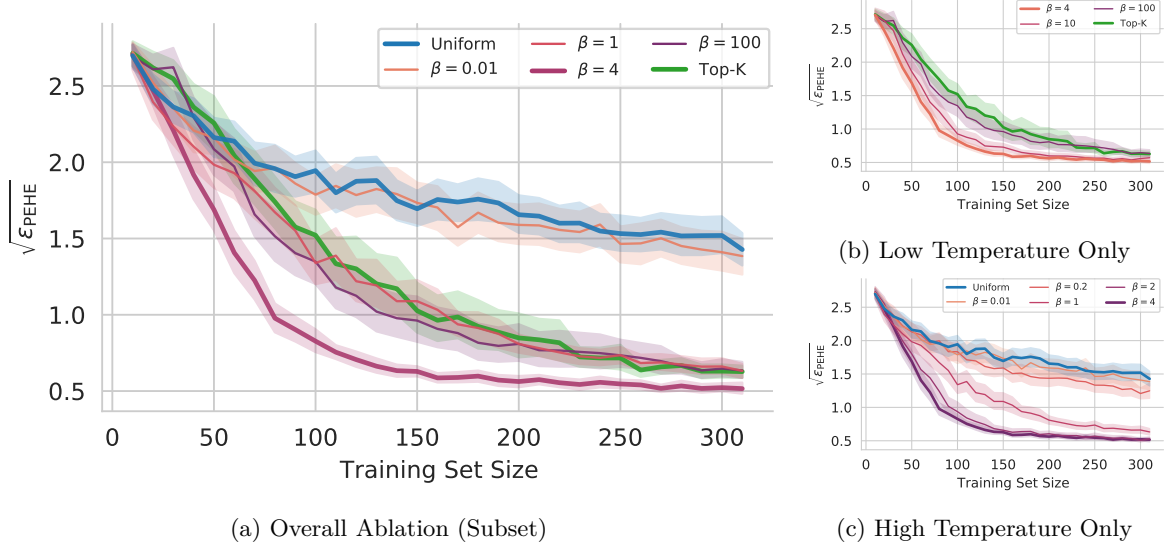(b) Low Temperature Only

(c) High Temperature Only

Figure 33: *CausalBALD: Synthetic Dataset.* (a) At a very high temperature ($\beta = 0.1$), PowerBALD behaves very much like random acquisition, and as the temperature decreases the performance of the acquistion function improves (lower $\sqrt{\epsilon_{\mathrm{PEHE}}}$). (b) Eventually, the performance reaches an inflection point ($\beta = 4.0$) and any further decrease in temperature results in the acquisition strategy performing more like top-$K$. We see that under the optimal temperature, power acquisition significantly outperforms both random acquisition and top-$K$ over a wide range of temperature settings.

We provide further $\beta$ ablations for CausalBALD on the entirely synthetic dataset which is used by Jesson et al. (2021). This demonstrates the ways in which $\beta$ interpolates between uniform and top-$K$ acquisition.
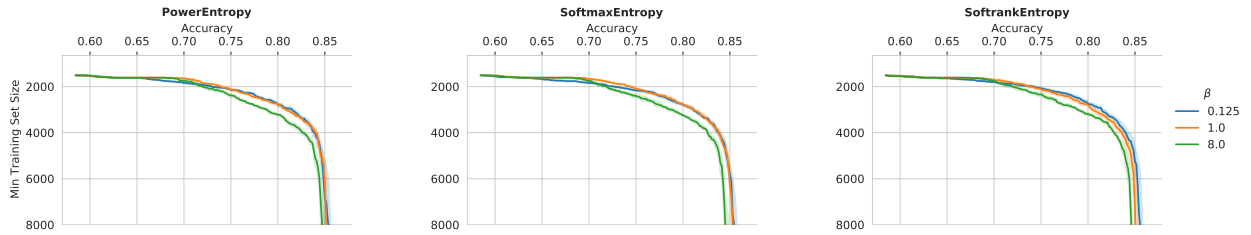
## E.3 CLINC-150



Figure 34: Performance CLINC-150: $\beta$ ablation for *Entropy.