

---

# Towards Understanding the Behaviors of Optimal Deep Active Learning Algorithms

---

Yilun Zhou\*

Sida Wang<sup>†</sup>

\*MIT CSAIL

Adithya Renduchintala<sup>†</sup>

Yashar Mehdad<sup>†</sup>

Xian Li<sup>†</sup>

Asish Ghoshal<sup>†</sup>

<sup>†</sup>Facebook AI

## Abstract

Active learning (AL) algorithms may achieve better performance with fewer data because the model guides the data selection process. While many algorithms have been proposed, there is little study on what the *optimal* AL algorithm looks like, which would help researchers understand where their models fall short and iterate on the design. In this paper, we present a simulated annealing algorithm to search for this optimal oracle and analyze it for several tasks. We present qualitative and quantitative insights into the behaviors of this oracle, comparing and contrasting them with those of various heuristics. Moreover, we are able to consistently improve the heuristics using one particular insight. We hope that our findings can better inform future active learning research. The code is available at <https://github.com/YilunZhou/optimal-active-learning>.

## 1 Introduction

Training deep models typically requires a large dataset, which limits its usefulness in domains where expensive expert annotation is required. Traditionally, active learning (AL), in which the model selects the data points to annotate and learn from, is often presented as a more sample efficient alternative to standard supervised learning (Settles, 2009). However the gain of AL with deep models is less consistent. For example, the best method seems to depend on the task in an unpredictable manner (Lowell et al., 2019).

Is active learning still useful in the era of deep learning? Although many issues are identified, it is not clear whether those problems only plague current AL methods or also apply to methods developed *in the*

*future*. Moreover, in existing literature, there lacks a comparison of proposed methods to the *oracle upper-bound*. Such a comparison is helpful for debugging machine learning models in many settings. For example, the classification confusion matrix may reveal the inability of the model to learn a particular class, and a comparison between a sub-optimal RL policy and the human expert play may indicate a lack of exploration. By contrast, without such an oracle reference in AL, it is extremely difficult, if at all possible, to pinpoint the inadequacy of an AL method and improve it.

In this paper, we propose a simulated annealing algorithm to search for the *optimal* AL strategy for a given base learner. With practical computational resources, this procedure is able to find an oracle that significantly outperforms existing heuristics (+7.53% over random orders, compared to +1.49% achieved by the best heuristics on average across three vision and language tasks), for models both trained from scratch and finetuned from pre-training, **definitively asserting the usefulness of a high-performing AL algorithm in most scenarios (Sec. 6.1)**. We also present the following insights into its behaviors.

While many papers do not explicitly state whether and how training stochasticity (e.g. model initialization or dropout) is controlled across iterations, we show that **training stochasticity tends to negatively affect the oracle performance (Sec. 6.2)**.

Previous work (Lowell et al., 2019) has found that for several heuristic methods, the actively acquired dataset does not transfer across different architectures. We observed a lesser extent of this phenomenon, but more importantly, **the oracle transfers better than heuristics (Sec. 6.3)**.

It may seem reasonable that a high-performing AL algorithm should exhibit a non-uniform sampling behavior (e.g. focusing more on harder to learn regions). However, **the oracle mostly preserves data distributional properties (Sec. 6.4)**.

Finally, using the previous insight, **heuristics can on average be improved by 2.95% with a simple distribution-matching regularization (Sec. 6.5)**.

## 2 Related Work

Active learning (Settles, 2009) has been studied for a long time. At the core of an active learning algorithm is the acquisition function, which generates or selects new data points to label at each iteration. Several different types of heuristics have been proposed, such as those based on uncertainty (Kapoor et al., 2007), disagreement (Houlsby et al., 2011), diversity (Xu et al., 2007), and expected error reduction (Roy and McCallum, 2001). In addition, several recent studies focused on meta-active learning, i.e. learning a good acquisition function on a source task and transfer it to a target task (Konyushkova et al., 2017; Fang et al., 2017; Woodward and Finn, 2017; Bachman et al., 2017; Contardo et al., 2017; Pang et al., 2018; Konyushkova et al., 2018; Liu et al., 2018; Vu et al., 2019).

With neural network models, active learning has been applied to computer vision (CV) (e.g. Gal et al., 2017; Kirsch et al., 2019; Sinha et al., 2019) and natural language processing (NLP) (e.g. Shen et al., 2018; Fang et al., 2017; Liu et al., 2018; Kasai et al., 2019) tasks. However, across different tasks, it appears that the relative performance of different methods can vary widely: a method can be the best on one task while struggling to even outperform the random baseline on another, with little explanation given. We compile a meta-analysis in App. A. Moreover, Lowell et al. (2019) found that the data acquisition order does not transfer well across architectures, a particularly important issue during deployment when it is expected that the acquired dataset will be used for future model development (i.e. datasets outliving models).

The closest work to ours is done by Koshorek et al. (2019), who studied the active learning limit. There are numerous differences. First, we analyze several CV and NLP tasks, while they focused on semantic role labeling (SRL) in NLP. Second, we explicitly account for training stochasticity, shown to be important in Sec. 6.2, but they ignored it. Third, our global simulated annealing search is able to find significant gaps between the upper limit and existing heuristics while their local beam search failed to find such a gap (though on a different task). We additionally show how to improve upon heuristics with our insights.

## 3 Problem Formulation

### 3.1 Active Learning Loop

Let the input and output space be denoted by  $\mathcal{X}$  and  $\mathcal{Y}$  respectively, where  $\mathcal{Y}$  is a finite set. Let  $\mathbb{P}_{XY}$  be the data distribution over  $\mathcal{X} \times \mathcal{Y}$ . We consider multi-class classification where a model  $m_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ , parameterized by  $\theta$ , maps an input to an output.

We study pool-based batch-mode active learning

where we assume access to an unlabeled data pool  $\mathcal{D}^U$  drawn from the marginal distribution  $\mathbb{P}_X$  over  $\mathcal{X}$ .

Starting with a labeled warm-start set  $\mathcal{D}_0^L \sim \mathbb{P}_{XY}$ , an AL loop builds a sequence of models  $(\theta_0, \dots, \theta_K)$  and a sequence of labeled datasets  $(\mathcal{D}_1^L, \dots, \mathcal{D}_K^L)$  in an interleaving manner. At the  $k$ -th iteration, for  $k = 0, \dots, K$ , a trainer  $\eta$  takes in  $\mathcal{D}_k^L$  and produces the model parameters  $\theta_k$  by minimizing some loss function on  $\mathcal{D}_k^L$ . In deep-learning the model training is typically stochastic (due to, for example, random initialization and drop-out masking) and  $\theta_k$  is a random variable. We assume that all such stochasticity are captured in  $\xi$  such that  $\theta_k = \eta(\mathcal{D}_k^L, \xi)$  is deterministic.

Using the trained model  $m_{\theta_k}$  and the current labeled set  $\mathcal{D}_k^L$ , an acquisition function  $\mathcal{A}$  builds the dataset  $\mathcal{D}_{k+1}^L$  by selecting a batch of  $B$  data points from the unlabeled pool,  $\Delta\mathcal{D}_{k+1} \subseteq \mathcal{D}^U$ , querying the annotator for labels, and adding them to  $\mathcal{D}_k^L$ ; i.e.  $\Delta\mathcal{D}_{k+1} = \mathcal{A}(\theta_k, \mathcal{D}_k^L)$ ;  $\mathcal{D}_{k+1}^L = \mathcal{D}_k^L \cup \Delta\mathcal{D}_{k+1}$ .<sup>1</sup> The model training and data acquisition loop repeats until we obtain  $\mathcal{D}_K^L$  by querying labels for  $KB$  data points, followed by training  $m_{\theta_K}$  for a final time.

### 3.2 Performance Curve

Given a specific draw of  $\xi$ , we measure the performance of an acquisition function  $\mathcal{A}$  by its *performance curve*  $\tau_{\mathcal{A}, \xi} : \{1, \dots, K\} \rightarrow [0, 1]$  defined below:

$$\tau_{\mathcal{A}, \xi}(k) = \mathbb{E}_{x, y \sim \mathbb{P}_{XY}} [e(m_{\theta_k}(x), y)],$$

where  $e : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, 1]$  is the evaluation metric (e.g. accuracy).  $\theta_k = \eta(\mathcal{D}_k^L, \xi)$ , and  $\mathcal{D}_k^L = \mathcal{D}_{k-1}^L \cup \mathcal{A}(\theta_{k-1}, \mathcal{D}_{k-1}^L)$ .

**Definition 1** (Anytime optimality). *An acquisition function  $\mathcal{A}$  is  $\xi$ -anytime optimal if it uniformly dominates every other acquisition function  $\mathcal{A}'$  as measured by  $\tau_{\cdot, \xi}$ ; i.e.,  $\tau_{\mathcal{A}, \xi}(k) \geq \tau_{\mathcal{A}', \xi}(k) \forall k \in \{1, \dots, K\}$  and  $\forall \mathcal{A}' \neq \mathcal{A}$ .*

**Proposition 1.** *There exist data distribution  $\mathbb{P}_{XY}$  and model class  $m_\theta$  for which an anytime optimal acquisition function does not exist.*

*Proof Sketch.* Fig. 1 shows a counter-example. In Fig. 1(a), we have an underlying distribution  $\mathbb{P}_{XY}$  shown as the colored background, and four points drawn from the distribution. If we learn a max-margin linear classifier from the four points, we can uncover the ground-truth decision boundary. If we acquire two data points, the optimal combination is shown in Fig. 1(b), resulting in a slightly wrong decision boundary. Choosing a different blue point would result in a much worse decision boundary (Fig. 1(c)). However,

<sup>1</sup>Some acquisition functions such as BALD are stochastic. For simplicity, we discuss the case of deterministic  $\mathcal{A}$  here. Extension of our claims to the stochastic case is straightforward and presented in App. B.

the optimal three-point acquisition (Fig. 1(d)), which leads to the ground-truth decision boundary, does *not* contain the previously optimal blue point in Fig. 1(b). Thus, there is no acquisition function simultaneously optimal at both two and three data points.  $\square$

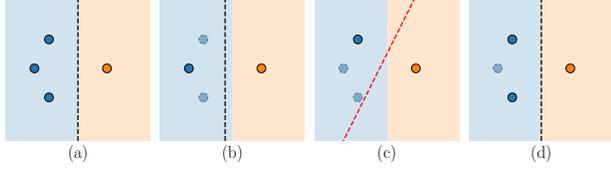


Figure 1: The counter example to prove Prop. 1.

**Remark.** In most AL papers, the authors demonstrate the superiority of their proposed method by showing its performance curve visually “above” those of the baselines. However, Prop. 1 shows that such an anytime optimal acquisition function may not exist.

### 3.3 $\xi$ -Quality

To make different acquisition functions comparable, we propose the following average performance metric, which we will refer to as  $\xi$ -quality:

$$q_\xi(\mathcal{A}) \stackrel{\text{def}}{=} \frac{1}{K} \sum_{k=1}^K \tau_{\mathcal{A},\xi}(k).$$

We refer to the acquisition strategy that maximizes this  $\xi$ -quality as  $\xi$ -optimal strategy, denoted as  $\mathcal{A}_\xi^*$ . Obviously, the  $\xi$ -anytime optimal acquisition function, if exists, will also be  $\xi$ -optimal.

While we could define the quality simply as the end point performance,  $\tau_{\mathcal{A},\xi}(K)$ , doing so fails to distinguish acquisition functions that lead to a fast initial performance rise from the rest, potentially incurring more annotation cost than necessary.

There are two interpretations to the above quality definition. First, it is the right Riemann sum approximation of the area under curve (AUC) of  $\tau_{\mathcal{A},\xi}$ , which correlates with intuitive notions of optimality for acquisition functions. Second, the un-averaged version of our definition can also be interpreted as the undiscounted cumulative reward of an acquisition *policy* over  $K$  steps where the per-step reward is the resulting model performance following data acquisition. Thus, the optimal policy implicitly trades off between immediate ( $k = 1$ ) and future ( $k \rightarrow K$ ) model performances.

### 3.4 Optimal Data Labeling Order

We index data points in  $\mathcal{D}^U$  by  $\mathcal{D}_i^U$  and use  $y(\mathcal{D}_i^U)$  to denote the label. Since an acquisition function generates a sequence of labeled datasets  $\mathcal{D}_1^L \subset \dots \subset \mathcal{D}_K^L$ , an acquisition function is equivalent to a permutation order  $\sigma$  of  $KB$  indices of  $\mathcal{D}^U$  with  $\mathcal{D}_k^L =$

$\mathcal{D}_0^L \cup \{(\mathcal{D}_{\sigma_i}^U, y(\mathcal{D}_{\sigma_i}^U))\}_{i \in [kB]}$ . Thus, the problem of searching for  $\mathcal{A}_\xi^*$  reduces to that of searching for the  $\xi$ -optimal order  $\sigma_\xi^*$ , which also relieves us from explicitly considering different forms of acquisition functions.

Moreover, a direct implication of above definitions is that optimal order depends on  $\xi$ ; i.e.  $q_\xi(\sigma_\xi^*) \geq q_{\xi'}(\sigma_{\xi'}^*)$ , for  $\xi \neq \xi'$ . As we experimentally demonstrate in Sec. 6.2, such a gap does exist. Since stochasticity in model training is completely controllable, an acquisition function that approaches optimal limit may need to explicitly take such randomness into account.

## 4 Search for the Optimal Order

There are two technical challenges in finding  $\sigma_\xi^*$ . First, we do not have access to the data distribution. Second, the permutation space of all orders is too large.

### 4.1 Validation Set Maximum

To solve the first problem, we assume access to a validation set  $\mathcal{D}^V \sim \mathbb{P}_{XY}$ . Since we are searching for the oracle model, there is no practical constraints on the size of  $\mathcal{D}^V$ . In addition, we assume access to an independently drawn test set  $\mathcal{D}^T \sim \mathbb{P}_{XY}$ . We define our optimal order estimator as

$$\hat{\sigma}_\xi = \arg \max_{\sigma} q_\xi(\sigma, \mathcal{D}^V), \quad (1)$$

and its quality on the test set  $q_\xi(\sigma, \mathcal{D}^T)$  serves as the an *unbiased* estimate of its generalization quality.

### 4.2 Simulated Annealing Search

Exhaustively searching over the space of all labeling orders is prohibitive as there are  $|\mathcal{D}^U|! / [(|\mathcal{D}^U| - BK)! \cdot (B!)^K]$  different orders to consider. In addition, we cannot use gradient-based optimization due to the discreteness of order space. Thus, we use simulated annealing (SA) (Kirkpatrick et al., 1983), as described in Alg. 1.

The search starts with a randomly initialized order  $\sigma^{(0)}$ . At time step  $t$ , it proposes a new order  $\sigma^{(p)}$  from  $\sigma^{(t-1)}$  with the following transition kernel (Fig. 2): with equal probabilities, it either swaps two data points from two different batches or replaces a data point in use by an unused one. It then evaluates the quality of the new order and accepts or rejects the proposal with probability depending on the quality difference. After searching for  $T_S$  steps, the algorithm greedily optimizes the best order for an additional  $T_G$  steps, and returns best order and quality in the end.

## 5 Experiment Overview

In addition to datasets  $\mathcal{D}^U, \mathcal{D}_0^L, \mathcal{D}^V$  and  $\mathcal{D}^T$  described in Sec. 3 and 4, since training neural networks on small

Task	Type	Dataset	$ \mathcal{D}^U ,  \mathcal{D}_0^L ,  \mathcal{D}^M ,  \mathcal{D}^V ,  \mathcal{D}^T $	$B, K$	Metric	Architecture	Heuristics
OC	class.	Fashion-MNIST	2000, 50, 150, 4000, 4000	25, 12	Acc	CNN	Max-Ent., (Batch)BALD
IC	class.	TOPv2 (alarm)	800, 40, 100, 4000, 4000	20, 8	F1	LSTM, CNN, AOE, RoBERTa	Max-Ent., BALD
NER	tagging	MIT Restaurant	1000, 50, 200, 3000, 3000	25, 10	F1	LSTM	(Norm-)Min-Conf., Longest

Table 1: Summary of experiment settings. Architecture details are in App. C.

**Algorithm 1:** Simulated Annealing (SA) Search

**Input:** SA steps  $T_S$ , greedy steps  $T_G$ , linear annealing factor  $\gamma$

$\sigma^{(0)} = \text{random-shuffle}([1, 2, \dots, |\mathcal{D}^U|]);$

$q^{(0)} = q_\xi(\sigma^{(0)});$

$\sigma^* = \sigma^{(0)}; q^* = q^{(0)};$

**for**  $t = 1, 2, \dots, T_S$  **do**

$\sigma^{(p)} = \text{propose}(\sigma^{(t-1)});$

$q^{(p)} = q_\xi(\sigma^{(p)}, \mathcal{D}^V);$

$u \sim \text{Unif}[0, 1];$

**if**  $u < \exp[\gamma t (q^{(p)} - q^{(t-1)})]$  **then**

$\sigma^{(t)} = \sigma^{(p)}; q^{(t)} = q^{(p)};$

**if**  $q^* < q^{(p)}$  **then**

$\sigma^* = \sigma^{(p)}; q^* = q^{(p)};$

**else**

$\sigma^{(t)} = \sigma^{(t-1)}; q^{(t)} = q^{(t-1)};$

**for**  $t = 1, 2, \dots, T_G$  **do**

$\sigma^{(p)} = \text{propose}(\sigma^*);$

$q^{(p)} = q_\xi(\sigma^{(p)}, \mathcal{D}^V);$

**if**  $q^{(p)} > q^*$  **then**

$\sigma^* = \sigma^{(p)}; q^* = q^{(p)};$

**return**  $\sigma^*, q^*$

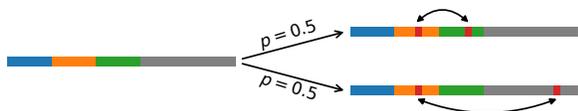


Figure 2: Illustration of the transition kernel. With equal probabilities, the kernel either swaps two data points from two different batches (sections of different colors) or replaces a data point in the first  $K$  batches with one outside (the gray section).

datasets is prone to overfitting, we introduced an additional model selection set  $\mathcal{D}^M$  to select the best model across 100 epochs and also trigger early stopping if the performance on this set is not improved for 20 consecutive epochs. In typical AL settings where data annotation is costly,  $\mathcal{D}^M$  is typically comparable in size to the warm-start set  $\mathcal{D}_0^L$ .

For each task, Tab. 1 describes the experiment-specific parameters. The training used the batch size equal to the acquisition batch size  $B$  and the Adam optimizer (Kingma and Ba, 2014). Following standard practice, the learning rate is  $10^{-3}$  except for RoBERTa fine-

tuning, which is  $10^{-4}$ . We searched for  $T_S = 25,000$  and  $T_G = 5,000$  steps with  $\gamma = 0.1$ . We found little improvement in the greedy stage.

The full dataset is shuffled and split into non-overlapping sets  $\mathcal{D}^U, \mathcal{D}_0^L, \mathcal{D}^M, \mathcal{D}^V, \mathcal{D}^T$ . Since the shuffling is *not* label-stratified, the empirical label distribution in any set (and  $\mathcal{D}_0^L$  in particular) may not be close to the actual distribution. We made this choice as it better reflects the actual deployment. For each task we considered commonly used architectures.

### 5.1 Object Classification (OC)

For object classification (OC), we used the Fashion-MNIST dataset (Xiao et al., 2017) with ten classes. Even though the dataset is label-balanced,  $|\mathcal{D}_0^L| = 50$  means that the initial warm-start set can be extremely imbalanced, posing additional challenge to uncertainty estimation used by many heuristics. We used a CNN architecture with two convolution layers followed by two fully connected layers. We compared against max-entropy, BALD (Gal et al., 2017), and BatchBALD (Kirsch et al., 2019) heuristics.

### 5.2 Intent Classification (IC)

For intent classification (IC), we used the Task-Oriented Parsing v2 (TOPv2) dataset (Chen et al., 2020), which consists of eight domains of human interaction with digital assistants, such as *weather* and *navigation*. We used the “alarm” domain for our experiments. In intent classification, given an user instruction such as “set an alarm at 8 am tomorrow”, the task is to predict the main intent of the sentence (**create-alarm**). There are seven intent classes: **create-alarm**, **get-alarm**, **delete-alarm**, **silence-alarm**, **snooze-alarm**, **update-alarm**, and **other**. The dataset is very imbalanced, with **create-alarm** accounting for over half of all examples. Hence, we used multi-class weighted F1 score as the evaluation metric.

Our main architecture of study is the bi-directional LSTM (BiLSTM) architecture with word embeddings initialized to GloVe (Pennington et al., 2014). To study the model transfer quality (Sec. 6.3), we also considered a CNN architecture, which uses 1D convolution layers to process a sentence represented as its word embedding sequence, an Average-of-Embedding (AOE) architecture, which uses a fully connected net-

work to process a sentence represented by the average embedding of the words, and finetuning from the pre-trained RoBERTa model (Liu et al., 2019). Detailed architectures are described in App. C. We considered max-entropy and BALD heuristics.

### 5.3 Named Entity Recognition (NER)

Named entity recognition (NER) is a structured prediction NLP task to predict a tag type for each word in the input sentence. We used the MIT Restaurant dataset (Liu et al., 2013), which consists of restaurant-related sentences tagged in Beginning-Inner-Outer (BIO) format. Tags include `amenity`, `price`, `location`, etc. For example, the tags for “*what restaurant near here serves pancakes at 6 am*” are `[0, 0, B-location, I-location, 0, B-dish, 0, B-hours, I-hours]`. The outer tag accounts for more than half of all tags, and tag-level weighted F1 score is used as the evaluation metric.

We used a BiLSTM encoder and an LSTM decoder following Shen et al. (2018) but without the CNN character encoder as an ablation study does not find it beneficial. We used teacher-forcing at training time and greedy decoding at prediction time. During active acquisition, we assumed that the annotation cost is the same for each sentence regardless of its length. We considered several heuristics. The min-confidence heuristic selects the sentence with the lowest log joint probability of the greedily decoded tag sequence (i.e. sum of log probability of each tag), which is divided by sentence length in its normalized version. The longest heuristic prioritizes longer sentences.

## 6 Results and Analyses

We compared the performance of the oracle order against existing heuristics. We called the oracle order and performance “optimal” even though they are estimates of the truly optimal ones due to our use of the validation set for estimating the quality and the simulated annealing search.

### 6.1 Optimal Quality Values and Curves

Tab. 2 presents the qualities of the optimal, heuristic and random orders. On all three tasks, there are large gaps between the heuristic and optimal orders. Specifically, the optimal orders are better than the best heuristics by 11.7%, 3.1%, and 3.5% for OC, IC, and NER respectively in terms of  $\xi$ -quality. For OC and IC, we repeated the training across five random seeds. The difference between the optimal and heuristic performances is significant at  $\alpha = 0.05$  for both tasks under a paired  $t$ -test, with  $p \leq 0.001$  and  $p = 0.003$  respectively. In addition, the optimal orders are better than the random orders by 7.53% on aver-

	OC	IC	NER
Optimal	<b>0.761</b>	<b>0.887</b>	<b>0.839</b>
Best Heuristic	0.682*	0.858	0.811
Random	0.698*	0.816	0.800

Table 2: Qualities of optimal, heuristic, and random orders on the three tasks. Individual heuristic performances are shown in App. D. \*No heuristics outperform the random baseline on object classification.

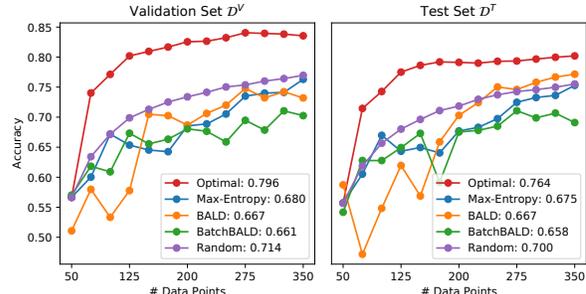


Figure 3: Performance curves of optimal, heuristic, and random orders for object classification.

age across three tasks, while the best heuristic orders only achieves an improvement of 1.49%.

Visually, Fig. 3 depicts the performance curve of the optimal order against heuristic and random baselines on both the validation ( $D^V$ ) and test ( $D^T$ ) set for OC. The optimal order significantly outperforms heuristic and random baselines, both numerically and visually, and we observed that the optimal order found using the validation set generalizes to the test set. Quality summary for each individual heuristic and performance curves for IC and NER tasks are presented in App. D and E. Even though there might not exist an anytime optimal acquisition function (Prop. 1), we did observe that the oracle order is uniformly better than heuristics on all three tasks.

### 6.2 Effect of Training Stochasticity

As previously stated, training stochasticity could negatively affect the optimal quality. To experimentally verify this, for  $\xi$  and  $\xi'$  we compared  $q_\xi(\hat{\sigma}_\xi, D^T)$  and  $q_\xi(\hat{\sigma}_{\xi'}, D^T)$  to study whether  $\xi'$ -optimal order is less optimal for  $\xi$ -training. Since we did not have dropout layers in our network, the only source of stochasticity comes from the random initialization.

For five different seeds in OC and IC,  $\xi, \xi' \in \{0, \dots, 4\}$ , Fig. 4 shows the pairwise quality  $q_\xi(\hat{\sigma}_{\xi'}, D^T)$ , where  $\xi'$  is the source seed on the column and  $\xi$  is the target seed on the row. The color on each cell represents the seed-mismatch quality gap  $q_\xi(\hat{\sigma}_\xi, D^T) - q_\xi(\hat{\sigma}_{\xi'}, D^T)$ , with a darker color indicating a larger gap.

Object Classification					Intent Classification					
0	0.764	0.748	0.753	0.754	0.748	0.887	0.880	0.884	0.883	0.884
1	0.755	0.760	0.752	0.753	0.751	0.872	0.882	0.880	0.879	0.879
2	0.757	0.746	0.764	0.760	0.750	0.883	0.880	0.891	0.882	0.889
3	0.749	0.746	0.750	0.762	0.749	0.883	0.881	0.883	0.884	0.882
4	0.742	0.729	0.743	0.742	0.756	0.876	0.879	0.882	0.883	0.890
	0	1	2	3	4	0	1	2	3	4
	Source Seed					Source Seed				

Figure 4: Training stochasticity negatively affects the optimal quality. The number in each cell represents  $q_{\xi}(\hat{\sigma}_{\xi'})$ , with  $\xi$  on the row and  $\xi'$  on the column. The color represents  $q_{\xi}(\hat{\sigma}_{\xi'}) - q_{\xi}(\hat{\sigma}_{\xi})$ , with darker colors indicating larger gaps.

The results suggest that in order to fully realize the potential of active learning, the acquisition function needs to specifically consider the particular model. However, high-quality uncertainty estimation are most commonly defined for a *model class* (e.g. the class of all random initialization in Deep Ensemble (Lakshminarayanan et al., 2017), the class of dropout masks in MC-Dropout (Gal and Ghahramani, 2016), and the class of Gaussian noises in Bayes-by-Backprop (Blundell et al., 2015)). A definition and algorithms for single-model uncertainty have the potential of further improving uncertainty-based heuristics. In addition, the result also implies that a purely diversity-based method relying on pre-defined embedding or distance metric could not be optimal.

### 6.3 Model Transfer Quality

Lowell et al. (2019) found that the heuristic acquisition order on one model sometimes does not transfer well to another model. However, it is not clear whether this is a general bane of active learning or just about specific heuristics. On the one hand, if the optimal order finds generally valuable examples early on, we would expect a different model architecture to also benefit from this order. On the other hand, if the optimal order exploits particularities of the model architecture, we would expect it to have worse transfer quality.

We replicated this study on the IC task, using four very different deep architectures: BiLSTM, CNN, AOE and RoBERTa described in Sec. 5. For a pair of source-target architectures ( $A_s, A_t$ ), we found both the optimal and best heuristic order on  $A_s$  and evaluate it on  $A_t$ . Then we compared its quality to that of a random order on  $A_t$ . The best heuristic is BALD for LSTM, CNN and RoBERTa, and max-entropy for AOE. Fig. 5 shows the quality of each order. The optimal order

		Optimal				Random	Heuristic				
Target Architecture	LSTM	0.887	0.869	0.877	0.845	0.817	0.860	0.850	0.852	0.838	1.0 0.0 -1.0
	CNN	0.781	0.853	0.811	0.764	0.754	0.781	0.799	0.794	0.756	
	AOE	0.847	0.855	0.879	0.814	0.798	0.845	0.848	0.848	0.818	
	RoBERTa	0.896	0.897	0.905	0.914	0.874	0.890	0.827	0.890	0.891	
		LSTM	CNN	AOE	RoBERTa		LSTM	CNN	AOE	RoBERTa	
		Source Architecture					Source Architecture				

Figure 5: Quality of optimal and max-entropy heuristic order across different architectures. Cell color represents difference to the random baseline.

always transfers better than the random order and, with the single exception of RoBERTa  $\rightarrow$  AOE, better than or equal to the heuristic order as well. This suggests that the optimal order tends to label model-agnostically valuable data points, and it is likely that better acquisition functions for one architecture can also help assemble high-quality datasets for others.

Perhaps surprisingly, we found that the optimal orders for different architectures actually share fewer data points than the heuristic orders, despite achieving higher transfer performance. Fig. 6 shows how the optimal (top) and heuristic (bottom) orders for different architectures relate to each other on the intent classification task. Specifically, each vertical bar is composed of 160 horizontal stripes, which represent the 160 data points acquired under the specified order, out of a pool of 800 data points (not shown). The first stripe from the top represents the first data point being acquired, followed by data points represented by the second, third, etc. stripes from the top. A line between the two orders connects the same data point in both orders and reflects the position of that data point in each of the two orders. The number of data points shared by both orders (i.e. number of connecting lines) is computed and shown.

The results suggest that there are an abundant amount of high quality data points that could lead to high AL performance, even though the optimal order for different architectures selects largely different ones. However, they are mostly missed by the heuristics. The heavy sharing among heuristics can be explained by the fact that they mostly over-sample sentences that are short in length or belong to the minority class, and thus deplete the pool of such examples.

### 6.4 Distributional Characteristics

In this set of experiments, we visualized the distributional characteristics of the acquired data points under the optimal and heuristic orders. We analyzed both the input and the output space.

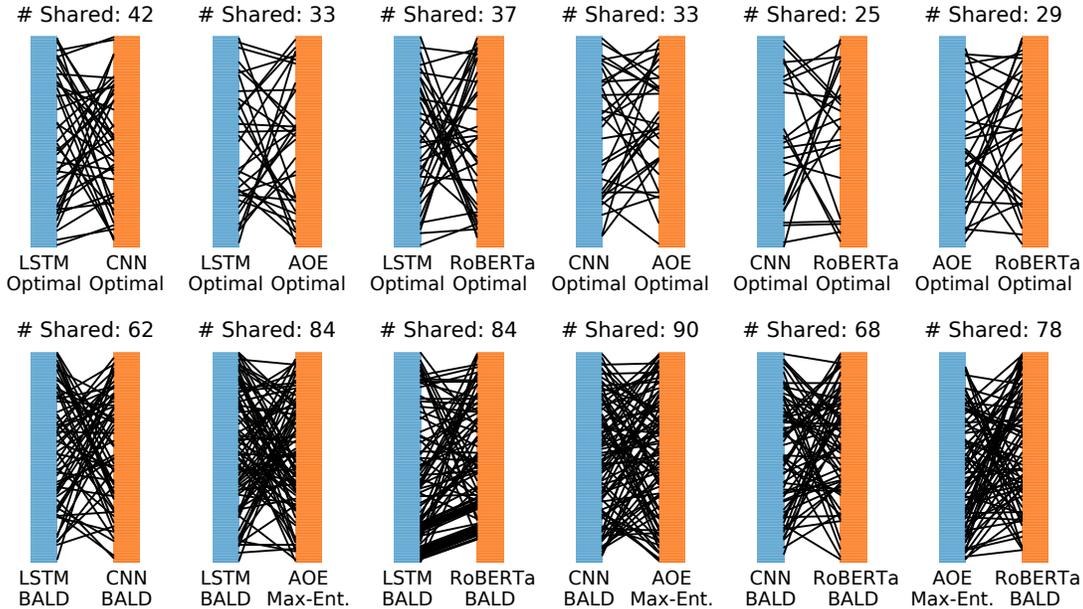


Figure 6: A visual comparison of optimal (top) and heuristic (bottom) orders, for every pair of architectures, showing the number of shared data points acquired under both architectures and their relative ranking.

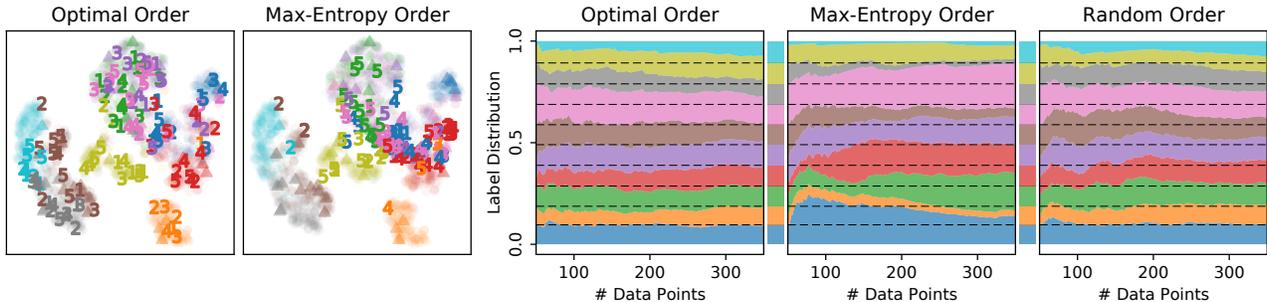


Figure 7: Distribution visualization for object classification. Left: first five batches numerically labeled in  $t$ -SNE embedding space, along with warm-start (triangle) and test (circle) samples. Right: label distribution w.r.t. labeled set size, with test-set distribution shown between plots and as dashed lines.

### 6.4.1 Object Classification

The left two panels of Fig. 7 visualizes the first five batches (125 data points in total) acquired by the optimal and max-entropy orders. For  $28 \times 28$  grayscale Fashion-MNIST images, we reduced their dimension to 100 using principal component analysis (PCA), then to 2 using  $t$ -distributed stochastic neighbor embedding ( $t$ -SNE) (Maaten and Hinton, 2008). Each acquired image, positioned at its 2-dimensional embedding and color-coded by label, is represented by a number referring to the batch index. Circle and triangle marks represent points in the test and warm-start sets.

The right three panels of Fig. 7 visualizes label distribution of the  $\mathcal{D}_k^L$  during acquisition. The horizontal axis represents the total number of data points (starting from 50 images in  $\mathcal{D}_0^L$ , each order acquires 300 from

the pool), while the relative length of each color represents the frequency of each label. The bars between plots show the “reference distribution” from the test set, which is balanced in this task.

For both input and output, the optimal order samples quite uniformly in the sample space, and is not easily distinguishable from random sampling. However, this order is distinctively non-random because it achieves a much higher quality than the random order (Tab. 2). By contrast, the max-entropy sampling is very imbalanced in both input and output space. BALD and BatchBALD exhibit similar behaviors (App. F).

### 6.4.2 Intent Classification

Fig. 8 presents the analogous analysis for IC, where we studied the sentence-length aspect of the input distri-

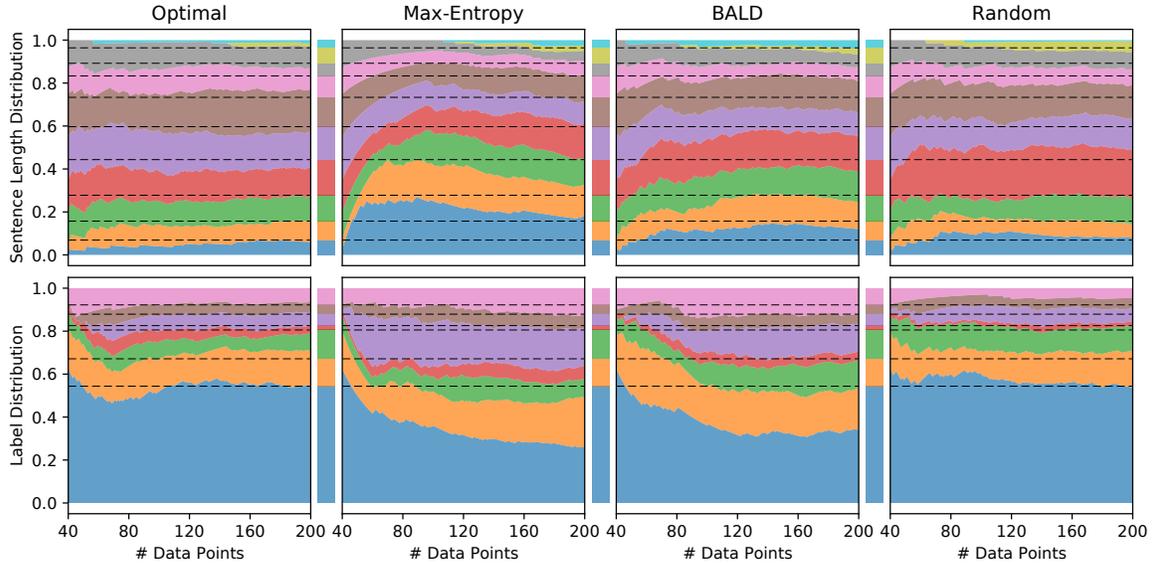


Figure 8: Sentence length and label distribution w.r.t. labeled set size for intent classification. For sentence lengths, the colors represent, from bottom to top: 1–3, 4, 5, 6, 7, 8, 9, 10, 11–12, and 13+. For labels, the colors represent, from bottom to top: `create-alarm`, `get-alarm`, `delete-alarm`, `other`, `silence-alarm`, `snooze-alarm`, and `update-alarm`.

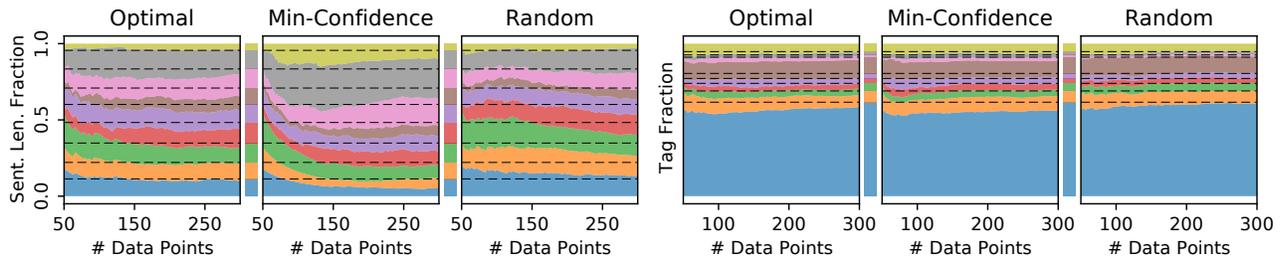


Figure 9: Sentence length and tag distribution w.r.t. labeled set size for NER. For sentence lengths, the colors represent, from bottom to top: 1–5, 6, 7, 8, 9, 10, 11–12, 13–16, and 17+. For labels, the colors represent, from bottom to top: `Outer`, `B/I-amenity`, `B/I-cuisine`, `B/I-dish`, `B/I-hours`, `B/I-location`, `B/I-price`, `B/I-rating`, and `B/I-restaurant-name`.

bution. Since there are few sentences with less than four or more than ten words, we grouped those sentences into 0-3, 11-12, and 13+ categories for ease of visualization. Both input and output distribution plots echo those in Fig. 7.

In particular, even though short sentences (blue and orange) are infrequent in the actual data distribution, both heuristics over-samples them, likely because shorter sentences provide less “evidence” for a prediction. In addition, they also significantly under-samples the `create-alarm` label (blue), likely because it is already well-represented in the warm-start set, and over-samples the `silence-alarm` label (purple), for the converse reason. The optimal order, again, does not exhibit any of these under- or over-sampling behaviors while achieving a much higher quality.

### 6.4.3 Named Entity Recognition

Fig. 9 presents the analysis for the NER task. We observed the same distribution-matching property, on both the input and the output side. By contrast, even if the acquisition function does not have the “freedom” to choose sentences with arbitrary tags, the min-confidence heuristic still over-samples the relatively infrequent `B/I-hours` tags (purple) in the beginning.

Notably, in NER the basic “supervision unit” is an individual word-tag pair. Thus, the model can effectively learn from more data by selecting longer sentences. This is indeed exploited by the min-confidence heuristic, where very long sentences are vastly over-sampled. However, the optimal order does not exhibit such length-seeking behavior, and instead matches the test set distribution quite closely.

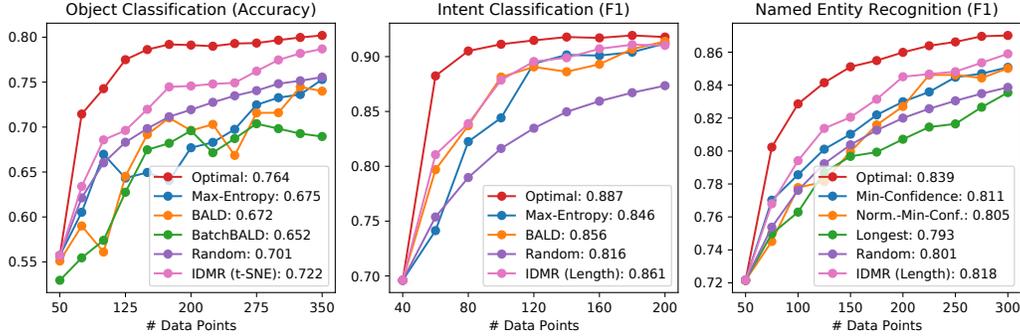


Figure 10: IDMR-augmented heuristics performance, improving upon the vanilla versions by 2.95% on average.

### 6.5 Distribution-Matching Regularization

Sec. 6.4 consistently suggests that unlike heuristics, the optimal order matches the data distribution meticulously. Can we improve heuristics with this insight?

---

#### Algorithm 2: Input Dist.-Matching Reg. (IDMR)

---

**Input:**  $\mathcal{A}(m_\theta, \mathcal{D}^L, \mathcal{D}^U)$  that returns the next data point in  $\mathcal{D}^U$  to label

$d_{\text{ref}} = \text{bin-distribution}(\mathcal{D}_{0,X}^L \cup \mathcal{D}_X^U \cup \mathcal{D}^M);$

$d_{\text{cur}} = \text{bin-distribution}(\mathcal{D}_X^L);$

$b^* = \arg \min_b (d_{\text{cur}} - d_{\text{ref}})_b;$

$\mathcal{D}_{b^*}^U = \{x \in \mathcal{D}_X^U : \text{bin}(x) = b^*\};$

**return**  $\mathcal{A}(m_\theta, \mathcal{D}^L, \mathcal{D}_{b^*}^U)$

---

Alg. 2 presents Input Distribution-Matching Regularization (IDMR), a meta-algorithm that augments any acquisition function to its input distribution-matching regularized version.  $\text{bin}(x)$  is a function that assigns each input data point  $x$  to one of finitely many bins by its characteristics.  $\text{bin-distribution}(\mathcal{D}_X)$  computes the empirical count distribution of the bin values for input part of  $\mathcal{D}_X$ . The IDMR algorithm compares the reference bin distribution on all available data (i.e.  $\mathcal{D}_{0,X}^L \cup \mathcal{D}_X^U \cup \mathcal{D}^M$ ) to the current bin distribution of  $\mathcal{D}_X^L$ , and uses the base acquisition function to select a data point in the highest deficit bin.

For OC, a  $K$ -means algorithm identified five clusters in the  $t$ -SNE space, which are used as the five bins. For IC and NER, sentences are binned by length. We used IDMR to augment the best performing heuristics. As shown in Fig. 10, on all tasks, the IDMR-augmented heuristic is better than all but the optimal order, with an average improvement of 2.95% over the its non-augmented counterpart. Across five random seeds for OC and IC, paired  $t$ -tests yield  $p = 0.28$  and  $0.017$  for improvement, both significant at  $\alpha = 0.05$ .

We also tried using a similar technique to match the output label distribution. However, it is complicated by the fact that the labels are not observed prior to acquisition and we observed mixed results. The algorithm and results are described in App. G.

## 7 Discussion and Conclusion

In this paper, we searched for and analyzed the optimal data acquisition order in active learning. Our findings are consistent across tasks and models. First, we confirmed the dependence of optimal orders on training stochasticity such as model initialization or dropout masking, which should be explicitly considered as AL methods approach the optimal limit.

Notably, we did not find any evidence that the optimal order needs to over-/under-sample in hard-/easy-to-learn regions. Given this observation, it is reasonable that the optimal order transfers across architectures better than heuristics and that heuristics can benefit from a distribution-matching regularization.

Indeed, most AL heuristics seek to optimize proxy objectives such as maximal informativeness or input space coverage. While intuitive, they have not been rigorously justified to correlate with any AL performance (e.g.  $\xi$ -quality). For example, even if a data point is maximally informative, it is possible that the information could not be fully absorbed during training and lead to performance improvement.

Moreover, in supervised learning, the fundamental assumption that the training and test data come from the same underlying distribution is the key to most guarantees of empirical risk minimization (ERM). Thus, we conjecture that the distribution-matching property arises from the nice behavior of ERM under this assumption. This also suggests that when the proposed algorithm is expected to violate this assumption, more careful analysis should be done on how such distribution shift may (adversely) affect the performance of models optimized under ERM. Developing theoretical guarantees of ERM under controlled distribution mismatch and/or formulations beyond ERM may also benefit active learning.

Finally, we focused on classification problems with a few classes. It would be interesting for future work to extend the analysis to other settings, such as large number of classes (e.g. CIFAR-100) or regression, and study whether the same behaviors hold or not.

## Acknowledgement

Yilun Zhou performed this work during an internship at Facebook AI. Correspondence should be sent to Yilun Zhou at [yilun@mit.edu](mailto:yilun@mit.edu). We thank the anonymous reviewers for their feedback.

## References

- Bachman, P., Sordoni, A., and Trischler, A. (2017). Learning algorithms for active learning. In *International Conference on Machine Learning (ICML)*, pages 301–310. PMLR.
- Blundell, C., Cornebise, J., Kavukcuoglu, K., and Wierstra, D. (2015). Weight uncertainty in neural network. In *International Conference on Machine Learning (ICML)*, pages 1613–1622. PMLR.
- Chen, X., Ghoshal, A., Mehdad, Y., Zettlemoyer, L., and Gupta, S. (2020). Low-resource domain adaptation for compositional task-oriented semantic parsing. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Contardo, G., Denoyer, L., and Artières, T. (2017). A meta-learning approach to one-step active-learning. In *International Workshop on Automatic Selection, Configuration and Composition of Machine Learning Algorithms*, volume 1998, pages 28–40. CEUR.
- Fang, M., Li, Y., and Cohn, T. (2017). Learning how to active learn: A deep reinforcement learning approach. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 595–605.
- Gal, Y. and Ghahramani, Z. (2016). Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *International Conference on Machine Learning (ICML)*, pages 1050–1059.
- Gal, Y., Islam, R., and Ghahramani, Z. (2017). Deep bayesian active learning with image data. In *International Conference on Machine Learning (ICML)*, pages 1183–1192. PMLR.
- Houlsby, N., Huszár, F., Ghahramani, Z., and Lengyel, M. (2011). Bayesian active learning for classification and preference learning. *arXiv preprint arXiv:1112.5745*.
- Kapoor, A., Grauman, K., Urtasun, R., and Darrell, T. (2007). Active learning with gaussian processes for object categorization. In *International Conference on Computer Vision (ICCV)*, pages 1–8. IEEE.
- Kasai, J., Qian, K., Gurajada, S., Li, Y., and Popa, L. (2019). Low-resource deep entity resolution with transfer and active learning. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 5851–5861.
- Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Kirkpatrick, S., Gelatt, C. D., and Vecchi, M. P. (1983). Optimization by simulated annealing. *science*, 220(4598):671–680.
- Kirsch, A., van Amersfoort, J., and Gal, Y. (2019). Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 7026–7037.
- Konyushkova, K., Sznitman, R., and Fua, P. (2017). Learning active learning from data. In *Advances in Neural Information Processing Systems (NIPS)*, pages 4225–4235.
- Konyushkova, K., Sznitman, R., and Fua, P. (2018). Discovering general-purpose active learning strategies. *arXiv preprint arXiv:1810.04114*.
- Koshorek, O., Stanovsky, G., Zhou, Y., Srikumar, V., and Berant, J. (2019). On the limits of learning to actively learn semantic representations. In *Conference on Computational Natural Language Learning (CoNLL)*, pages 452–462.
- Lakshminarayanan, B., Pritzel, A., and Blundell, C. (2017). Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems (NIPS)*, pages 6402–6413.
- Liu, J., Pasupat, P., Cyphers, S., and Glass, J. (2013). ASGAR: A portable architecture for multilingual dialogue systems. In *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8386–8390. IEEE.
- Liu, M., Buntine, W., and Haffari, G. (2018). Learning how to actively learn: A deep imitation learning approach. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1874–1883.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. (2019). Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Lowell, D., Lipton, Z. C., and Wallace, B. C. (2019). Practical obstacles to deploying active learning. In *Conference on Empirical Methods in Natural Language Processing and International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 21–30.
- Maaten, L. v. d. and Hinton, G. (2008). Visualizing data using *t*-SNE. *Journal of Machine Learning Research (JMLR)*, 9(Nov):2579–2605.

- Pang, K., Dong, M., Wu, Y., and Hospedales, T. (2018). Meta-learning transferable active learning policies by deep reinforcement learning. *arXiv preprint arXiv:1806.04798*.
- Pennington, J., Socher, R., and Manning, C. D. (2014). Glove: Global vectors for word representation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543.
- Roy, N. and McCallum, A. (2001). Toward optimal active learning through monte carlo estimation of error reduction. *International Conference on Machine Learning (ICML)*, pages 441–448.
- Settles, B. (2009). Active learning literature survey. Technical report, University of Wisconsin-Madison Department of Computer Sciences.
- Shen, Y., Yun, H., Lipton, Z. C., Kronrod, Y., and Anandkumar, A. (2018). Deep active learning for named entity recognition. In *International Conference on Learning Representation (ICLR)*.
- Sinha, S., Ebrahimi, S., and Darrell, T. (2019). Variational adversarial active learning. In *International Conference on Computer Vision (ICCV)*, pages 5972–5981.
- Vu, T., Liu, M., Phung, D., and Haffari, G. (2019). Learning how to active learn by dreaming. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 4091–4101.
- Woodward, M. and Finn, C. (2017). Active one-shot learning. *arXiv preprint arXiv:1702.06559*.
- Xiao, H., Rasul, K., and Vollgraf, R. (2017). Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms.
- Xu, Z., Akella, R., and Zhang, Y. (2007). Incorporating diversity and density in active learning for relevance feedback. In *European Conference on Information Retrieval (ECIR)*, pages 246–257. Springer.