

FLEXIBLE ACTIVE LEARNING OF PDE TRAJECTORIES

Anonymous authors

Paper under double-blind review

ABSTRACT

Accurately solving partial differential equations (PDEs) is critical for understanding complex scientific and engineering phenomena, yet traditional numerical solvers are computationally expensive. Surrogate models offer a more efficient alternative, but their development is hindered by the cost of generating sufficient training data from numerical solvers. In this paper, we present a novel framework for active learning (AL) in PDE surrogate modeling that reduces this cost. Unlike the existing AL methods for PDEs that always acquire entire PDE trajectories, our approach strategically generates only the most important time steps with the numerical solver, while employing the surrogate model to approximate the remaining steps. This dramatically reduces the cost incurred by each trajectory and thus allows the active learning algorithm to try out a more diverse set of trajectories given the same budget. To accommodate this novel framework, we develop an acquisition function that estimates the utility of a set of time steps by approximating its resulting variance reduction. We demonstrate the effectiveness of our method on several benchmark PDEs, including the Heat equation, Korteweg–De Vries equation, Kuramoto–Sivashinsky equation, and the incompressible Navier-Stokes equation. Extensive experiments validate that our approach outperforms existing methods, offering a cost-efficient solution to surrogate modeling for PDEs.

1 INTRODUCTION

In many scientific and engineering applications, accurately solving partial differential equations (PDEs) in the form of trajectories of states evolving over time is essential for understanding complex phenomena (Holton & Hakim, 2013; Atkins et al., 2023; Murray, 2007; Wilmott et al., 1995). The traditional approach involves running numerical solvers, which provide accurate solutions but are computationally costly, taking several hours, days or even weeks to run depending on the complexity of the problem (Cleaver et al., 2016; Cowan et al., 2001). As a result, there is significant interest in developing surrogate models (Greydanus et al., 2019; Bar-Sinai et al., 2019; Sanchez-Gonzalez et al., 2020; Li et al., 2020; Brandstetter et al., 2022b; Lippe et al., 2024) that can approximate the solutions more efficiently. Surrogate models are obtained by solving regression tasks on some “ground truth” data. The ground truth data for PDEs are generated by numerical solvers, which are costly compared to those of standard regression problems. As a result, the expense of data acquisition presents a major bottleneck in the development of surrogate models for PDEs.

Active Learning (AL, Chernoff, 1959; MacKay, 1992; Settles, 2009) can address this challenge by adaptively acquiring the most informative inputs, effectively reducing the amount of ground-truth data required to obtain a high-quality surrogate model. However, there is a general lack of research in AL for regression tasks (Wu, 2018; Holzmüller et al., 2023), let alone PDEs. Existing studies on AL for PDEs have predominantly dealt with univariate outputs such as energy (Pestourie et al., 2020; Pickering et al., 2022), or predictions at a single, fixed time point (Bajracharya et al., 2024; Wu et al., 2023b). To our surprise, the only work directly addressing AL for prediction of trajectories is that of Musekamp et al. (2024). In this work, the surrogate model is set as an autoregressive model that predicts the evolved state of a PDE at time $t + \Delta t$ given a state at an arbitrary time point t , and is trained on data acquired by existing regression-based AL methods (Holzmüller et al., 2023). Specifically, at each round of acquisition, the AL method chooses initial conditions from which *entire* trajectories are acquired. However, we argue that querying all the states in a trajectory is not **sample-efficient**, especially for autoregressive surrogate models.

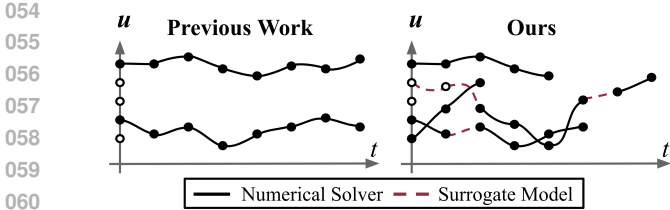


Figure 1: Conceptual illustration of our framework for data acquisition. Each dot represents a PDE state, and a path connecting two dots represents a time step of a simulation. Black solid lines are obtained with a numerical solver and red dotted lines with a surrogate model.

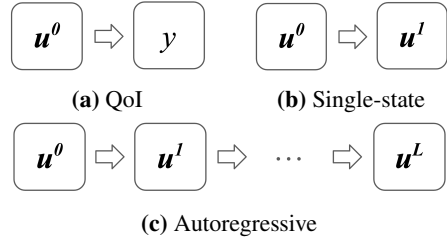


Figure 2: Task settings assumed by previous works in active learning of PDEs.

Acquiring entire trajectories is inefficient mainly for two reasons. First, states within a trajectory are often strongly correlated, undermining their diversity or the joint information gain (Houlsby et al., 2011; Kirsch et al., 2019). We validate this assertion in Appendix C.1 through principal component analysis. Secondly, even if they are not strongly correlated, it can be the case that only certain time steps of a trajectory are the most informative due to the dynamics of the PDE. In both cases, noting that the main cost is in running the numerical solver, it would be ideal to selectively acquire only the most important time steps with the numerical solver, for a fraction of the cost of acquiring the entire trajectory. However, this is usually impossible without querying all the time steps that come earlier.

In this paper, we propose a novel, flexible framework for data acquisition that circumvents the constraint of having to query all time steps in a trajectory, along with an AL strategy that leverages this flexibility. Our method combines both a numerical solver and a surrogate model to acquire data along a trajectory with reduced cost. Specifically, it selects which time steps along a trajectory to query from the solver, while using the surrogate model to approximate the remaining steps. We also develop a novel acquisition function that guides our AL strategy in choosing which time steps to query to the numerical solver in each trajectory.

Overall, our framework, equipped with the novel AL strategy, significantly improves surrogate model performance over previous methods. We validate our approach through extensive experiments on benchmarks, including the Heat equation, Korteweg–De Vries equation, Kuramoto–Sivashinsky equation, and the Navier–Stokes equation. Additionally, we analyze the behavior of our AL method, providing insights into the factors that contribute to its effectiveness.

2 BACKGROUND

2.1 PRELIMINARIES

We consider PDEs with one time dimension $t \in [0, T]$ and possibly multiple spatial dimensions $\mathbf{x} = [x_1, x_2, \dots, x_D] \in \mathbb{X}$ where \mathbb{X} is the spatial domain such as the unit interval. These can be written in the form

$$\partial_t \mathbf{u} = F(t, \mathbf{x}, \mathbf{u}, \partial_{\mathbf{x}} \mathbf{u}, \partial_{\mathbf{x}\mathbf{x}} \mathbf{u}, \dots), \tag{1}$$

where $\mathbf{u} : [0, T] \times \mathbb{X} \rightarrow \mathbb{R}^n$ is a solution to the PDE. We are also given a specific boundary condition and a fixed time interval Δt . If the PDE is well-posed (Evans, 1988), there exists, for each $t_0 \in \mathbb{R}$, an evolution operator G_{t_0} which maps an initial condition $\mathbf{u}^0 := \mathbf{u}(t_0, \cdot)$ to the solution $\mathbf{u}^1 := \mathbf{u}(t_0 + \Delta t, \cdot)$. For simplicity, we only consider time-independent PDEs, for which the evolution operator G_t is the same for all t , say G . Iterating over G multiple times, we can obtain a trajectory $(\mathbf{u}^i)_{i=1}^L$ of length L , where $\mathbf{u}^i := G^{(i)}[\mathbf{u}^0]$ with $G^{(i)}$ being the i -th iterate of G . Although a numerical solver G_{solver} is only an approximation to G , we shall not distinguish between the two for the remainder of this paper.

There are three primary tasks in active learning for PDEs, each depending on the type of surrogate model being trained. The first task, *univariate Quantity of Interest (QoI) prediction*, focuses on

learning a model to directly predict a scalar QoI, denoted as y , from an initial condition \mathbf{u}^0 . The second task, *single-state prediction*, involves learning a model to predict a single state transition from \mathbf{u}^0 to \mathbf{u}^1 over a fixed time interval Δt . The third task, *autoregressive trajectory prediction*, aims to approximate the ground truth evolution operator G using a surrogate model to predict the entire time evolution of the states. Fig. 2 provides a visual comparison of the three tasks. In this paper, we focus on the autoregressive trajectory prediction task.

We train a neural surrogate model \hat{G} with input-output pairs $(\mathbf{u}^{i-1}, \mathbf{u}^i)$ from the numerical solver G . Active learning builds a high quality training dataset by adaptively selecting informative inputs to be fed into the solver G . Prior work (Musekamp et al., 2024) operates on the the framework where initial conditions \mathbf{u}^0 are selected from a pool \mathcal{P} , from which full trajectories of length L are obtained. For instance, Query-by-Committee (QbC, Seung et al., 1992) queries initial conditions \mathbf{u}^0 that maximize the predictive uncertainty estimated from a committee of M models,

$$a_{\text{QbC}}(\mathbf{u}^0) = \frac{1}{M} \sum_{m=1}^M \sum_{i=1}^L \|\hat{\mathbf{u}}_m^i - \bar{\mathbf{u}}^i\|_2^2 \quad (2)$$

where $\hat{\mathbf{u}}_m^i$ is the prediction of the i^{th} state from the m^{th} surrogate model in the committee and $\bar{\mathbf{u}}^i := \frac{1}{M} \sum_{m=1}^M \hat{\mathbf{u}}^i$ is the mean prediction from the committee.

2.2 PROBLEM SETTING

Our ultimate objective is to obtain a surrogate model \hat{G} that approximates the expensive numerical solver G with low error

$$\frac{1}{N_{\text{test}}} \sum_{j=1}^{N_{\text{test}}} \text{err} \left((G^{(i)}[\mathbf{u}_j^0])_{i=1}^L, (\hat{G}^{(i)}[\mathbf{u}_j^0])_{i=1}^L \right) \quad (3)$$

where $\text{err}(\cdot, \cdot)$ is an error metric. Obtaining the surrogate model requires sampling training data from the numerical solver, which incurs a nontrivial cost. AL aims to improve sample efficiency by sampling only the most important data. In particular, AL utilizes the current surrogate model \hat{G} , or a committee of surrogate models $\{\hat{G}_m\}_{m=1}^M$, to inform its choice. After acquiring the data chosen by AL, we retrain the surrogate \hat{G} with the expanded training dataset.

We assume that there exists a pool \mathcal{P} of initial conditions \mathbf{u}^0 . At each round of AL, we train a committee of M surrogate models $\{\hat{G}_m\}_{m=1}^M$ with the training dataset collected from G up to that round. We then use this committee to select a batch of inputs to be queried to the solver G and add them to the training dataset. The cost at each round, defined as the number of inputs queried to the numerical solver G , is limited to a certain budget B . Our aim is to achieve low errors at each round, so an AL strategy would ideally acquire data with cost as close to or equal to the budget (Li et al., 2022a).

3 FLEXIBLE ACTIVE LEARNING FOR PDES

3.1 FRAMEWORK OF DATA ACQUISITION

We present our method, FLEXAL, which operates under a framework of data acquisition that is much more sample efficient than previous works. Algorithm 1 provides an overview of our framework. We start with a surrogate model \hat{G} , or a committee of surrogate models $\{\hat{G}_m\}_{m=1}^M$, trained with the initial dataset \mathcal{D} . At every round of AL, we choose an initial condition \mathbf{u}^0 from the pool \mathcal{P} , similar to the existing AL methods for PDE trajectories. However, while existing methods acquire the entire trajectory starting from the chosen initial condition \mathbf{u}^0 (Musekamp et al., 2024), our method acquires a partial trajectory. Specifically, we select a subset of time steps to simulate from \mathbf{u}^0 , rather than acquiring the full trajectory. The rationale behind this approach is that, given a fixed budget, acquiring as many trajectories as possible—albeit partially—from different initial conditions is often more beneficial than fully acquiring fewer trajectories. This strategy enables more efficient exploration of the data space and improves the overall sample efficiency of the framework.

Algorithm 1 Overview of Flexible Active Learning (FLEXAL)

Require: Pool \mathcal{P} of initial conditions, budget B per round, number of rounds R , numerical solver G , trajectory length L , initial training dataset \mathcal{D}

Ensure: Trained surrogate model \hat{G}

- 1: Train \hat{G} on \mathcal{D}
- 2: **for** round = 1 to R **do**
- 3: cost $\leftarrow 0$
- 4: **while** cost $< B$ **do**
- 5: Choose initial condition \mathbf{u}^0 from \mathcal{P} ▷ Section 3.3
- 6: $\mathcal{P} \leftarrow \mathcal{P} \setminus \{\mathbf{u}^0\}$
- 7: Choose sampling pattern $S = (b_1, \dots, b_L)$ ▷ Sections 3.2 and 3.3
- 8: $\hat{\mathbf{u}}^0 \leftarrow \mathbf{u}^0$
- 9: **for** $i = 1$ to L **do**
- 10: $\hat{\mathbf{u}}^i \leftarrow \begin{cases} G[\hat{\mathbf{u}}^{i-1}] & \text{if } b_i = \text{true} \\ \hat{G}[\hat{\mathbf{u}}^{i-1}] & \text{if } b_i = \text{false} \end{cases}$
- 11: **if** $b_i = \text{true}$ **then**
- 12: $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\hat{\mathbf{u}}^{i-1}, \hat{\mathbf{u}}^i)\}$
- 13: cost \leftarrow cost + 1
- 14: **end if**
- 15: **end for**
- 16: **end while**
- 17: Train \hat{G} on \mathcal{D}
- 18: **end for**

More specifically, for a given initial condition \mathbf{u}^0 , we define a boolean sequence of length L , $S = (b_1, \dots, b_L)$, which we refer to as the *sampling pattern*. For example, S could be (true, false, ..., true). The sampling pattern specifies that data will be acquired only at time steps corresponding to true values while skipping those marked false.

After selecting the sampling pattern S , the next step is to acquire the PDE trajectory. While acquiring a full trajectory is straightforward using a numerical solver G , obtaining a partial trajectory corresponding to S can be tricky. We want to run the solver G only for the time steps specified by S (those with true patterns), but the solver requires the skipped time steps (those with false patterns) as intermediate inputs. If we just run the solver for all time steps for this reason, we wouldn't be saving any cost. To address this, we use a simple heuristic: for the skipped time steps, we replace the simulation with predictions from the *surrogate model* (we use the average surrogate $\hat{G} = \frac{1}{M} \sum_{m=1}^M \hat{G}_m$ when we have a committee). That is, starting with $\hat{\mathbf{u}}^0 = \mathbf{u}^0$, we iterate over $1 \leq i \leq L$:

$$\hat{\mathbf{u}}^i = \begin{cases} G[\hat{\mathbf{u}}^{i-1}] & \text{if } b_i = \text{true} \\ \hat{G}[\hat{\mathbf{u}}^{i-1}] & \text{if } b_i = \text{false}. \end{cases} \quad (4)$$

We add to our dataset \mathcal{D} only the input-output pairs obtained with the solver G , namely $(\hat{\mathbf{u}}^{i-1}, \hat{\mathbf{u}}^i)$ with $b_i = \text{true}$.

In comparison to full trajectory acquisition, which requires L executions of the numerical solver, our strategy invokes the numerical solver $\|S\| := \sum_{i=1}^L \mathbb{1}[b_i = \text{true}]$ times and utilizes the surrogate model $L - \|S\|$ times. Since the surrogate model is significantly cheaper to evaluate than the numerical solver, this approach substantially reduces the cost of acquisition, enabling us to explore more initial conditions within the same budget. In fact, as discussed in Section 2.2, we define the acquisition cost precisely as $\|S\|$. We repeat expanding our training dataset with new initial conditions and sampling patterns until the cost incurred in the current round reaches a budget B . At the end of each round, we retrain the surrogate \hat{G} with the expanded training dataset \mathcal{D} .

Previous methods listed in Musekamp et al. (2024) can be considered a special case of ours where the sampling pattern S is always full of true entries. Our framework is therefore a strict generalization of previous works. In the remainder of this section, we describe how FLEXAL adaptively chooses initial conditions \mathbf{u}^0 and sampling patterns S .

3.2 ACQUISITION FUNCTION

To adaptively select the sampling pattern S with the initial condition \mathbf{u}^0 , we propose a novel acquisition function $a(\mathbf{u}^0, S)$ that assesses the utility of S . Given a committee $\{\hat{G}_m\}_{m=1}^M$, consider (\hat{G}_a, \hat{G}_b) for some $a, b \in [M] := \{1, \dots, M\}$ with $a \neq b$. We define the utility of the sampling pattern S for the pair (\hat{G}_a, \hat{G}_b) as the resulting *variance reduction* in the pair’s rolled-out trajectories. Specifically, let $\hat{\mathbf{u}}_a$ and $\hat{\mathbf{u}}_b$ be the trajectories estimated by \hat{G}_a and \hat{G}_b , starting from \mathbf{u}^0 . Next, we obtain a rollout using Eq. 4 with our sampling pattern S and surrogate model \hat{G}_b , where \hat{G}_a serves as a stand-in for the ground-truth solver G . We denote the resulting trajectory as $\hat{\mathbf{u}}_{b,S,a}$. The variance reduction is defined as

$$R(a, b, S) := \sum_{i=1}^L (\|\hat{\mathbf{u}}_a^i - \hat{\mathbf{u}}_b^i\|^2 - \|\hat{\mathbf{u}}_a^i - \hat{\mathbf{u}}_{b,S,a}^i\|^2). \quad (5)$$

The sampling pattern S that maximizes $R(a, b, S)$ is the one where the current models \hat{G}_a and \hat{G}_b disagree the most, and acquiring data from S effectively reduces this discrepancy. Our acquisition function is defined as the average variance reduction between all the distinct pairs in the committee:

$$a(\mathbf{u}^0, S) = \frac{1}{M(M-1)} \sum_{a,b \in [M], a \neq b} R(a, b, S). \quad (6)$$

We observe that our acquisition function simplifies to QbC in Eq. 2 when S acquires all the time steps, differing only by a constant factor of two. This occurs because, in that case, $\hat{\mathbf{u}}_{b,S,a} = \hat{\mathbf{u}}_a$, which makes the second term in the summand of Eq. 2 vanish. Consequently, we can interpret our acquisition function as a generalization of QbC that accommodates for the selection of time steps.

As an additional sanity check, consider the scenario where S does not sample any time steps. In this situation, $\hat{\mathbf{u}}_{b,S,a} = \hat{\mathbf{u}}_b$, leading the two terms in the summand to cancel each other out, resulting in zero variance reduction. Since acquiring no data should yield zero utility, we confirm that our acquisition function behaves as expected in this limiting case. [Appendix D.1 further details the precise motivation behind the design of our acquisition function.](#)

3.3 BATCH ACQUISITION ALGORITHM

With the acquisition function defined above, we present a batch acquisition algorithm given a pool \mathcal{P} of initial conditions. We define the cost of a batch $\{(\mathbf{u}_j^0, S_j)\}_{j=1}^N$ as the total number of queries to the solver, $\sum_{j=1}^N \|S_j\|$. A standard objective is to maximize $\sum_{j=1}^N a(\mathbf{u}_j^0, S_j)$ under a budget constraint $\sum_j \|S_j\| \leq B$, to which there is a known approximate solution (Salkin & De Kluuyer, 1975) that greedily maximizes the cost-weighted acquisition function $a^*(\mathbf{u}^0, S) = a(\mathbf{u}^0, S) / \|S\|$ until the total cost exceeds the budget B . However, this method faces two problems. First, it’s questionable whether the sum $\sum_{j=1}^N a(\mathbf{u}_j^0, S_j)$ of individual acquisition values is actually a good representative for the utility of a batch. In fact, numerous works report that picking instances that maximize individual acquisition values can severely underperform compared to methods that take into account the interactions between those instances (Kirsch et al., 2019; Ash et al., 2019). The problem is chiefly attributed to the lack of diversity and representativeness (Wu, 2018) caused by oversampling of small, high value regions (Smith et al., 2023). Secondly, we are actually searching over the *product* pool of the sampling pattern S and the pool \mathcal{P} of initial conditions, whose size is on the order of $O(2^L |\mathcal{P}|)$. Both terms impose significant computational burden on optimizing the cost-weighted objective $a^*(\mathbf{u}^0, S)$.

We therefore propose FLEXAL as an add-on to existing AL methods that acquire full trajectories. Specifically, a full-trajectory AL method \mathcal{A} , which we call a *base* method, first selects an initial condition \mathbf{u}^0 . [Musekamp et al. \(2024\) introduces several possibilities for such a method, including QbC \(Seung et al., 1992\), Largest Cluster Maximum Distance \(LCMD, Holzmüller et al., 2023\), Core-Set \(Sener & Savarese, 2017\), and stochastic batch active learning \(SBAL, Kirsch et al., 2023\).](#) We then optimize the cost-weighted acquisition function $a(\mathbf{u}^0, S)$ over the sampling pattern S while holding \mathbf{u}^0 fixed, and add the pair (\mathbf{u}^0, S) to the current batch. We iterate this two-stage process until the cost of the batch reaches our budget limit. Additionally, if the cost ever exceeds

the budget after adding a pair, we truncate the sampling pattern so that the cost is exactly equal to the budget. By using FLEXAL as an add-on, the diversity and representativeness promoted by base AL methods (Holzmüller et al., 2023; Kirsch et al., 2023; Musekamp et al., 2024) are upheld, and the size of the optimization space for FLEXAL is reduced to $O(2^L)$. The problem remains, however, that $O(2^L)$ is a prohibitively large space for optimization. We therefore use a simple greedy algorithm for searching S . In the greedy algorithm, we start by initializing S with all entries set to true. At each step of the greedy algorithm, we propose a neighboring pattern S' by applying a bit-flip mutation, where each bit of S is flipped with a probability of ϵ . The proposal is accepted only if the acquisition value $a^*(u^0, S')$ is higher than the current value $a^*(u^0, S)$. This process of proposal and acceptance/rejection is repeated T times. We use $T = 100$ and $\epsilon = 0.1$ throughout our experiments. A more concise summary of the batch acquisition algorithm is given in Appendix D.2. The algorithmic complexity of batch acquisition is discussed in Section 5.7.

4 RELATED WORK

AL for PDEs. The works by Pestourie et al. (2020); Pickering et al. (2022); Gajjar et al. (2022) apply active learning to problems involving PDEs, but their tasks are limited to predicting QoI, such as the maximum value of an evolved state. Li et al. (2024); Wu et al. (2023b) apply their AL methods to single-state prediction. Bajracharya et al. (2024) explores the use of active learning in tasks of predicting steady states of PDEs, which can be seen as predicting single states at $t \rightarrow \infty$. Finally, Musekamp et al. (2024) experiments with active learning in predicting PDE trajectories with autoregressive models.

Active selection of time points. While our work is the first to propose time step selection in active learning (AL) for PDEs, the concept of selecting time points has been explored in other contexts. For example, in physics-informed neural networks (PINNs), active selection of collocation points for training has been widely studied (Arthurs & King, 2021; Gao & Wang, 2023; Mao & Meng, 2023; Wu et al., 2023a; Turinici, 2024). “Labels” for PINNs, or the residual loss, can be calculated directly at any time point using closed form equations. There are also methods in Bayesian experimental design (BED) that choose observation times that maximize information gain about parameters of interest (Singh et al., 2005; Cook et al., 2008). In those works, a trajectory is already “there”, but the cost is attributed to the act of observing a time point. In contrast, in our setting, we cannot directly acquire a time point, because there is a cost in the simulation of the trajectory.

Multi-fidelity AL. Closely related to our work is multi-fidelity active learning Li et al. (2022b); Wu et al. (2023b); Hernandez-Garcia et al. (2023); Li et al. (2024), where outputs are acquired at varying fidelity levels for each input, with associated costs inherent to each fidelity. In our context, the task of actively selecting a sampling pattern for a given initial condition can be seen as a fidelity selection problem, where acquiring all time steps corresponds to the highest fidelity but also incurs the highest cost.

5 EXPERIMENTS

5.1 BASELINE AL METHODS

To compare with our method, we experiment with AL for full trajectory sampling introduced in Musekamp et al. (2024). **Random** sampling from the pool set is the simplest method. **QbC** (Seung et al., 1992) is a simple active learning algorithm that selects points according to maximum disagreement among members of a committee. **LCMD** (Holzmüller et al., 2023) is an AL algorithm that uses a feature map. We concatenate the last hidden layer activations of committee members at all time steps of a trajectory, and sketch the concatenated features to a dimension of 512 using a random projection. Kirsch et al. (2023) proposes **SBAL**, which randomly samples data points x with a probability distribution proportional to its temperature-scaled acquisition value $p(x) \propto a(x)^m$. We use the acquisition function of QbC with temperature $m = 1$. We leave out Core-Set (Sener & Savarese, 2017) because it generally underperforms compared to the above methods, according to both Holzmüller et al. (2023) and Musekamp et al. (2024).

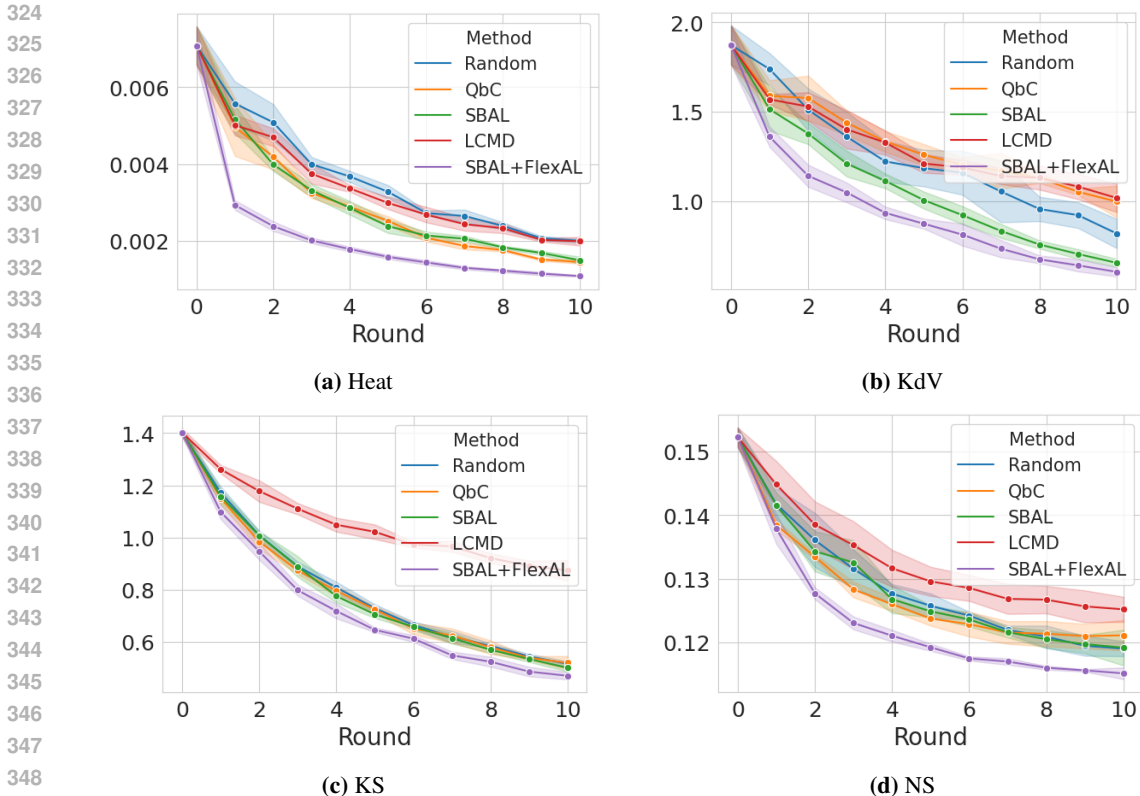


Figure 3: RMSE of AL strategies, measured across 10 rounds of acquisition. Each round incurs constant cost of data acquisition, namely the budget B .

5.2 TARGET PDES

We evaluate our method on a range of PDEs. The first is the **Heat** equation in one spatial dimension. Next, we test the nonlinear Korteweg–De Vries (**KdV**) equation, which is known for exhibiting solitary wave pulses with weak interactions (Zabusky & Kruskal, 1965). We then apply our method to the Kuramoto–Sivashinsky (**KS**) equation, another nonlinear PDE in one dimension, notable for its chaotic dynamics. Lastly, we consider the vorticity form of the incompressible Navier–Stokes equation (**NS**) in two spatial dimensions. All equations are solved with periodic boundary conditions. Additional details are in Appendix A.1.

5.3 SURROGATE MODELS

We use a Fourier Neural Operator (FNO, Li et al., 2020) to model the evolution operator G . In particular, we train it to predict the differences between states in adjacent time steps, following Musekamp et al. (2024). All models have four hidden layers. We use 16, 256, 128, and 32 modes for Heat, KdV, KS, and NS equations, respectively. We also normalize the data according to the initial dataset’s mean and standard deviation over all temporal and spatial dimensions. We use teacher-forcing to train the FNOs, meaning that it’s simply trained on ground truth input-output pairs from the solver G without backpropagating through two or more time steps. All models were trained with Adam (Kingma, 2014) for 100 epochs, using a learning rate of 10^{-3} , a batch size of 32, and a cosine annealing scheduler (Loshchilov & Hutter, 2016).

5.4 RESULTS

We compare between the four baselines introduced in Section 5.1, and our method combined with SBAL (SBAL+FLEXAL). The pool set has 10,000 initial conditions, and we always start with an initial dataset of 32 fully sampled trajectories. The initial conditions in the test set are sampled from

Table 1: Log RMSE of baseline methods and SBAL+FLEXAL, averaged across 10 rounds of acquisition

	Random	QbC	LCMD	SBAL	SBAL+FLEXAL
Heat	-5.688 ± 0.021	-5.924 ± 0.025	-5.741 ± 0.024	-5.901 ± 0.017	-6.304 ± 0.015
KdV	0.191 ± 0.058	0.266 ± 0.027	0.256 ± 0.030	0.030 ± 0.029	-0.088 ± 0.040
KS	-0.258 ± 0.003	-0.268 ± 0.003	0.046 ± 0.013	-0.275 ± 0.014	-0.349 ± 0.003
NS	-2.050 ± 0.011	-2.057 ± 0.009	-2.018 ± 0.017	-2.052 ± 0.009	-2.092 ± 0.003

Table 2: Log RMSE of FLEXAL averaged across 10 rounds, and their improvement over base methods. Δ refers to the improvements from baselines. Negative Δ indicates better performance of FLEXAL.

	Random		QbC		LCMD	
	+FLEXAL	Δ	+FLEXAL	Δ	+FLEXAL	Δ
Heat	-6.193 ± 0.021	-0.505 ± 0.030	-6.195 ± 0.015	-0.271 ± 0.029	-6.131 ± 0.024	-0.390 ± 0.034
KdV	-0.067 ± 0.054	-0.258 ± 0.079	0.134 ± 0.035	-0.132 ± 0.044	0.286 ± 0.034	0.030 ± 0.045
KS	-0.335 ± 0.012	-0.077 ± 0.012	-0.331 ± 0.013	-0.063 ± 0.013	-0.138 ± 0.016	-0.184 ± 0.021
NS	-2.080 ± 0.003	-0.030 ± 0.011	-2.079 ± 0.008	-0.022 ± 0.012	-2.050 ± 0.007	-0.032 ± 0.018

the same distribution as those in the pool set. An ensemble size of $M = 2$ is used, as it has been shown to be sufficient for good AL performance (Pickering et al., 2022; Musekamp et al., 2024). We perform 10 rounds of acquisition, and the budget of each round is set to $B = 8 \times L$ where L is the length of a trajectory. This means that full trajectory algorithms sample 8 trajectories per round. We report their RMSE, defined in Appendix A.2. Reports of other metrics are provided in Appendix B. Fig. 3 shows plots of the committee’s mean RMSE across the 10 rounds of acquisition, and Table 1 summarizes the results with mean logarithmic RMSEs, where a mean is taken over all 10 rounds. We can observe from the plots that SBAL+FLEXAL outperforms other AL baselines in a robust manner. Most notably, it improves the surrogate models on both the KS and NS equations, where no other baseline improves significantly over random sampling. On NS, SBAL+FLEXAL achieves an RMSE below 0.12 at the fifth round, which is only achieved by the best baseline at the tenth round. FLEXAL has effectively halved the cost of acquisition required to obtain this accuracy. All experiments were conducted on 8 NVIDIA GeForce RTX 2080 Ti GPUs, and the results are averages from 5 seed values.

5.5 OTHER BASE METHODS

We also report the mean log RMSE of FLEXAL when combined with the three other base methods, in Table 2. We find that the performance always improves over the base method with the addition of FLEXAL, except for the case of LCMD on KdV where the difference is negligible. Fig. 4 shows three plots of RMSE on the NS equation, where each contains a base method and its combination with FLEXAL. We find that the discrepancy between the two are noticeably larger for base methods that did not perform robustly when used alone. For instance, QbC tends to perform worse than Random in the later rounds, which is also where the discrepancy between QbC and QbC+FLEXAL becomes more noticeable. Also, LCMD performs the worst when used alone, and also creates the largest improvement when FLEXAL is added. We can infer that adding FLEXAL has the effect of swinging back to some loss curve, and that this effect is stronger for base methods that deviate more from it. We do note, however, that the loss curves of FLEXAL are distinct for different base methods.

5.6 RANDOM BERNOULLI SAMPLING OF TIME STEPS

We plot in Fig. 5 the distribution of time steps that our method chooses. The plot clearly shows the general tendency of FLEXAL to acquire the early time steps, with an occasional selection of the later time steps. The distributions still show clear differences between tasks, such as in their average number of time steps per trajectory or the frequency of later time steps. These suggests that FLEXAL is choosing time steps in an adaptive manner that’s different for each task at hand.

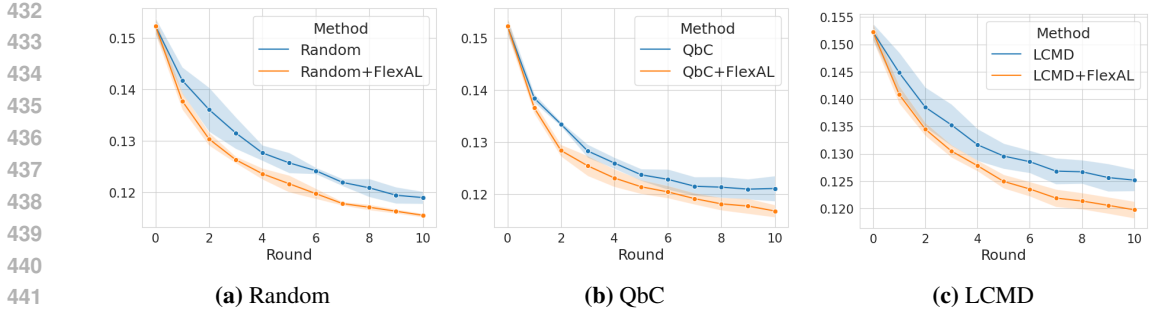


Figure 4: RMSE of base methods with and without FLEXAL on NS, measured across 10 rounds of acquisition

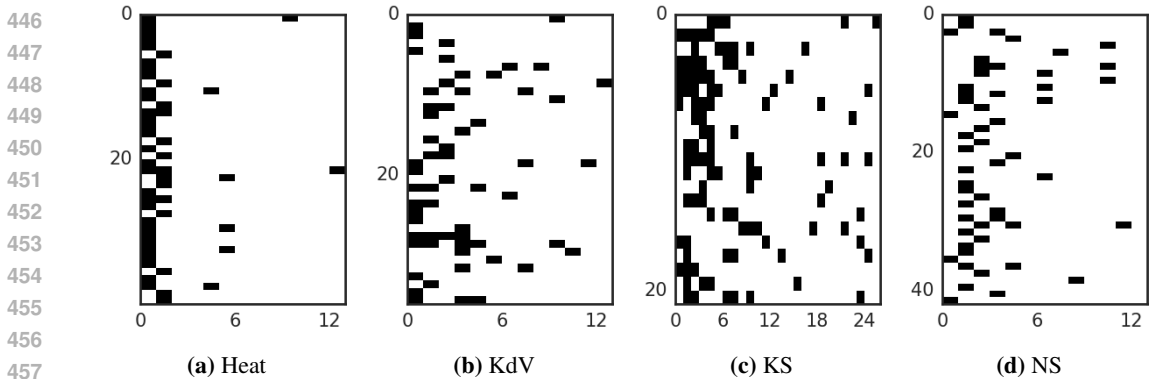


Figure 5: Timesteps chosen by SBAL+FLEXAL. Each row corresponds to an acquired trajectory, where the black cells indicate the selected time steps. Half of all trajectories acquired in the first rounds of active learning are shown.

We then ask ourselves: what if we perform random selection, for instance with a probability p , for every time step? We call this method Bernoulli sampling, or $\text{Ber}(p)$, where each entry of S is true with probability p . Table 3 summarizes the performance of $\text{Ber}(p)$ for $p = 1/16, 1/8, 1/4$, and $1/2$. Results show that FLEXAL outperforms Bernoulli sampling, except for the case of KS where $\text{Ber}(1/16)$ serves as a strong alternative. In general, Bernoulli sampling improves over the base method SBAL, but it can also severely underperform at certain values of p , such as for KdV. Still, for each PDE, there exists a value of p at which Bernoulli sampling provides an advantage over the base method SBAL. These observations show altogether that sparse sampling of time steps itself has an inherent advantage over full-trajectory sampling, and that FLEXAL amplifies this gain by adaptively choosing not only the frequency of the time steps to acquire, but also their locations. We report the full results in Appendix C.3, along with a variant of Bernoulli sampling that enforces acquiring consecutive initial time steps.

5.7 ALGORITHMIC COMPLEXITY OF FLEXAL

The time complexity of computing our acquisition function for a single instance of (\mathbf{u}^0, S) is $O(M^2L)$. Since we optimize the acquisition function with T steps, and we can acquire at most B initial conditions, the time complexity of our batch acquisition algorithm is $O(M^2LBT)$ in the worst case. We can parallelize the optimization of multiple S_j 's to a certain extent using graphics processing unit (GPU), which can significantly alleviate the burden of B . We can further reduce the cost by at most a factor of M with FLEXAL MF described in Appendix A.4. Yet another alternative is to decrease the number of greedy optimization steps T from 100 to 10, which reduces the cost by a factor of 10. We call this variant FLEXAL 10. The wall-clock time of each baseline method and FLEXAL is summarized in Table 4. The performance of SBAL with FLEXAL and its two variants are summarized in Appendix C.4, as well as the wall-clock times on all equations. Note

Table 3: Log RMSE with Bernoulli sampling averaged across 10 rounds of acquisition

	SBAL	+FLEXAL	+Ber(1/16)	+Ber(1/8)	+Ber(1/4)	+Ber(1/2)
Heat	-5.901 ± 0.017	-6.304 ± 0.015	-6.093 ± 0.018	-6.071 ± 0.020	-6.057 ± 0.026	-6.010 ± 0.035
KdV	0.030 ± 0.029	-0.088 ± 0.040	0.053 ± 0.014	0.049 ± 0.014	0.018 ± 0.024	-0.064 ± 0.031
KS	-0.275 ± 0.014	-0.349 ± 0.003	-0.365 ± 0.008	-0.359 ± 0.006	-0.346 ± 0.008	-0.324 ± 0.007
NS	-2.052 ± 0.009	-2.092 ± 0.003	-2.088 ± 0.005	-2.081 ± 0.008	-2.079 ± 0.007	-2.075 ± 0.009

Table 4: Wall-clock time of each procedure during batch selection in NS. Measured with a single NVIDIA GeForce RTX 2080 Ti GPU. **Note that these are not the costs of data acquisition, but the computational cost of batch selection algorithms.**

	Random	QbC	LCMD	SBAL	+FLEXAL	+FLEXAL MF	+FLEXAL 10
Time taken (seconds)	0.1 ± 0.1	45.5 ± 0.2	72.2 ± 1.4	45.1 ± 0.2	92.2 ± 2.6	55.9 ± 0.4	10.5 ± 1.2

that FLEXAL 10 incurs only a fraction of computational cost over the baseline methods, while still achieving a significant boost in performance over its base method. **After all, the increased computational cost of the selection process is negligible in practical settings because the cost of data acquisition usually far exceeds the cost of selection. In fact, without running the numerical solver in batch mode, obtaining data for a single round in the KdV experiment takes around 20 minutes, which is far greater than any of the costs incurred by the selection algorithms. Moreover, increasing the pool size increases the runtime of base methods, but doesn't incur any additional runtime on FLEXAL.**

6 CONCLUSION

In this paper, we presented a novel framework for active learning in surrogate modeling of partial differential equation (PDE) trajectories, significantly reducing the cost of data acquisition while maintaining or improving model accuracy. By selectively querying only a subset of time steps in a PDE trajectory, our method FLEXAL enables the acquisition of informative data at a fraction of the cost of acquiring entire trajectories. We introduced a new acquisition function that estimates the utility of a set of time steps based on variance reduction, effectively guiding the selection process in an adaptive manner. Through extensive experiments on benchmark PDEs, including the Heat equation, Korteweg–De Vries equation, Kuramoto–Sivashinsky equation, and incompressible Navier-Stokes equation, we demonstrated that our approach consistently outperforms existing AL methods, providing a more cost-efficient and accurate solution for PDE surrogate modeling.

Our results show that FLEXAL can significantly enhance surrogate modeling in PDEs, particularly in scenarios where the numerical solver is computationally expensive. We further showed that the success of FLEXAL is driven by its ability to prioritize both diverse and informative time steps. Moving forward, this framework could be extended to more complex systems and integrated with other machine learning techniques, providing broader applicability in scientific and engineering simulations. Future work may also explore alternative acquisition functions and applications to simulations outside the domain of PDEs.

540 REPRODUCIBILITY STATEMENT.

541
542 We present detailed description of our algorithm in [Section 3.3](#). Details regarding the algorithm’s
543 hyperparameters, model architecture, training, active learning procedure, and data generation are
544 provided in [Section 5](#) and [Appendix A](#).

546 ETHICS STATEMENT.

547
548 We propose a new method that improves the cost efficiency of acquiring data for building a surrogate
549 model of PDE trajectories. Although our approach doesn’t have a direct positive or negative impact
550 in ethical or societal aspects, it accelerates the process of building a surrogate model for an arbi-
551 trary PDE. This could be used for good, such as medical simulations, environmental modeling, and
552 optimizing engineering designs, potentially leading to advancements in healthcare, sustainability,
553 and technological innovation. However, like many technologies, this method could also be misused
554 in domains where rapid simulations could have harmful consequences, such as the development of
555 hazardous materials. Therefore, researchers and practitioners should apply these methods with con-
556 sideration of their broader societal implications, aiming to ensure that the benefits of the technology
557 are used responsibly and ethically.

559 REFERENCES

- 560
561 Christopher J Arthurs and Andrew P King. Active training of physics-informed neural networks to
562 aggregate and interpolate parametric solutions to the navier-stokes equations. *Journal of Computa-*
563 *tional Physics*, 438:110364, 2021. [6](#)
- 564 Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal.
565 Deep batch active learning by diverse, uncertain gradient lower bounds. *arXiv preprint*
566 *arXiv:1906.03671*, 2019. [5](#)
- 567 Peter William Atkins, Julio De Paula, and James Keeler. *Atkins’ physical chemistry*. Oxford univer-
568 sity press, 2023. [1](#)
- 569 Pradeep Bajracharya, Javier Quetzalcóatl Toledo-Marín, Geoffrey Fox, Shantenu Jha, and Linwei
570 Wang. Feasibility study on active learning of smart surrogates for scientific simulations. *arXiv*
571 *preprint arXiv:2407.07674*, 2024. [1](#), [6](#)
- 572 Yohai Bar-Sinai, Stephan Hoyer, Jason Hickey, and Michael P Brenner. Learning data-driven dis-
573 cretizations for partial differential equations. *Proceedings of the National Academy of Sciences*,
574 116(31):15344–15349, 2019. [1](#)
- 575 Johannes Brandstetter, Max Welling, and Daniel E Worrall. Lie point symmetry data augmentation
576 for neural pde solvers. In *International Conference on Machine Learning*, pp. 2241–2256. PMLR,
577 2022a. [15](#), [16](#)
- 578 Johannes Brandstetter, Daniel Worrall, and Max Welling. Message passing neural pde solvers. *arXiv*
579 *preprint arXiv:2202.03376*, 2022b. [1](#), [20](#)
- 580 Eric P Chassignet, Harley E Hurlburt, Ole Martin Smedstad, George R Halliwell, Patrick J Hogan,
581 Alan J Wallcraft, Remy Baraille, and Rainer Bleck. The hycom (hybrid coordinate ocean model)
582 data assimilative system. *Journal of Marine Systems*, 65(1-4):60–83, 2007. [16](#)
- 583 Herman Chernoff. Sequential design of experiments. *The Annals of Mathematical Statistics*, 30(3):
584 755–770, 1959. [1](#)
- 585 Timothy A Cleaver, Alex J Gutman, Christopher L Martin, Mark F Reeder, and Raymond R Hill.
586 Using design of experiments methods for applied computational fluid dynamics: A case study.
587 *Quality Engineering*, 28(3):280–292, 2016. [1](#)
- 588 Alex R Cook, Gavin J Gibson, and Christopher A Gilligan. Optimal observation times in experi-
589 mental epidemic processes. *Biometrics*, 64(3):860–868, 2008. [6](#)

- 594 Timothy J Cowan, Andrew S Arena Jr, and Kajal K Gupta. Accelerating computational fluid dy-
595 namics based aeroelastic predictions using system identification. *Journal of Aircraft*, 38(1):81–87,
596 2001. 1
- 597 WD Evans. Partial differential equations, 1988. 2
- 599 Aarshvi Gajjar, Chinmay Hegde, and Christopher P Musco. Provable active learning of neural
600 networks for parametric pdes. In *The Symbiosis of Deep Learning and Differential Equations II*,
601 2022. 6
- 602 Wenhan Gao and Chunmei Wang. Active learning based sampling for high-dimensional nonlinear
603 partial differential equations. *Journal of Computational Physics*, 475:111848, 2023. 6
- 604 Michael A Gelbart, Jasper Snoek, and Ryan P Adams. Bayesian optimization with unknown con-
605 straints. *arXiv preprint arXiv:1403.5607*, 2014. 17
- 606 Somdatta Goswami, Katiana Kontolati, Michael D Shields, and George Em Karniadakis. Deep
607 transfer operator learning for partial differential equations under conditional shift. *Nature Ma-
608 chine Intelligence*, 4(12):1155–1164, 2022. 17
- 609 Samuel Greydanus, Misko Dzamba, and Jason Yosinski. Hamiltonian neural networks. *Advances
610 in neural information processing systems*, 32, 2019. 1
- 611 Alex Hernandez-Garcia, Nikita Saxena, Moksh Jain, Cheng-Hao Liu, and Yoshua Bengio. Multi-
612 fidelity active learning with gflownets. *arXiv preprint arXiv:2306.11715*, 2023. 6
- 613 José Miguel Hernández-Lobato, Michael Gelbart, Matthew Hoffman, Ryan Adams, and Zoubin
614 Ghahramani. Predictive entropy search for bayesian optimization with unknown constraints. In
615 *International conference on machine learning*, pp. 1699–1707. PMLR, 2015. 17
- 616 James R Holton and Gregory J Hakim. *An introduction to dynamic meteorology*, volume 88. Aca-
617 demic press, 2013. 1
- 618 David Holzmüller, Viktor Zaverkin, Johannes Kästner, and Ingo Steinwart. A framework and bench-
619 mark for deep batch active learning for regression. *Journal of Machine Learning Research*, 24
620 (164):1–81, 2023. 1, 5, 6, 17, 19
- 621 Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. Bayesian active learning for
622 classification and preference learning. *arXiv preprint arXiv:1112.5745*, 2011. 2
- 623 Nicolas Jarrin, Sofiane Benhamadouche, Dominique Laurence, and Robert Prosser. A synthetic-
624 eddy-method for generating inflow conditions for large-eddy simulations. *International Journal
625 of Heat and Fluid Flow*, 27(4):585–593, 2006. 16
- 626 John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger,
627 Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate
628 protein structure prediction with alphafold. *nature*, 596(7873):583–589, 2021. 16
- 629 E Kalnay. *Atmospheric Modeling, Data Assimilation and Predictability*, volume 341. Cambridge
630 University Press, 2003. 16
- 631 Aly-Khan Kassam and Lloyd N Trefethen. Fourth-order time-stepping for stiff pdes. *SIAM Journal
632 on Scientific Computing*, 26(4):1214–1233, 2005. 16
- 633 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,
634 2014. 7
- 635 Andreas Kirsch, Joost Van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch
636 acquisition for deep bayesian active learning. *Advances in neural information processing systems*,
637 32, 2019. 2, 5
- 638 Andreas Kirsch, Sebastian Farquhar, Parmida Atighehchian, Andrew Jesson, Frédéric Branchaud-
639 Charron, and Yarin Gal. Stochastic batch acquisition: A simple baseline for deep active learning.
640 *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. 5, 6

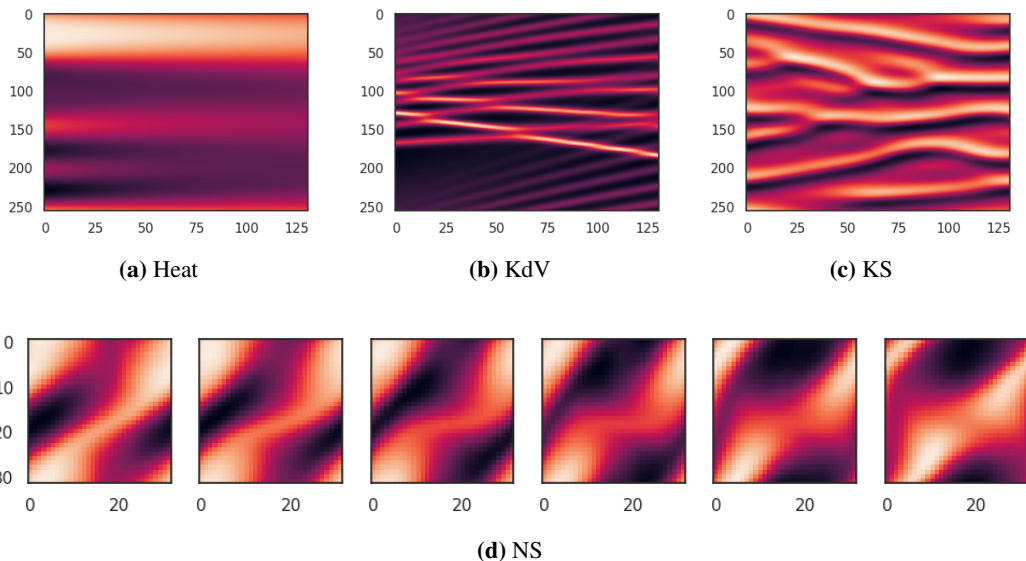
- 648 Matt J Kusner, Brooks Paige, and José Miguel Hernández-Lobato. Grammar variational autoen-
649 coder. In *International conference on machine learning*, pp. 1945–1954. PMLR, 2017. 16
650
- 651 Shibo Li, Jeff M Phillips, Xin Yu, Robert Kirby, and Shandian Zhe. Batch multi-fidelity active
652 learning with budget constraints. *Advances in Neural Information Processing Systems*, 35:995–
653 1007, 2022a. 3
- 654 Shibo Li, Zheng Wang, Robert M. Kirby, and Shandian Zhe. Deep multi-fidelity active learning of
655 high-dimensional outputs. In *International Conference on Artificial Intelligence and Statistics*,
656 *AISTATS 2022, 28-30 March 2022, Virtual Event, 2022b*. 6, 21
657
- 658 Shibo Li, Xin Yu, Wei Xing, Robert Kirby, Akil Narayan, and Shandian Zhe. Multi-resolution
659 active learning of fourier neural operators. In *International Conference on Artificial Intelligence*
660 *and Statistics*, pp. 2440–2448. PMLR, 2024. 6
- 661 Zongyi Li, Nikola Kovachki, Kamyar Aizzadenesheli, Burigede Liu, Kaushik Bhattacharya, An-
662 drew Stuart, and Anima Anandkumar. Fourier neural operator for parametric partial differential
663 equations. *arXiv preprint arXiv:2010.08895*, 2020. 1, 7, 16
664
- 665 Phillip Lippe, Bas Veeling, Paris Perdikaris, Richard Turner, and Johannes Brandstetter. Pde-refiner:
666 Achieving accurate long rollouts with neural pde solvers. *Advances in Neural Information Pro-*
667 *cessing Systems*, 36, 2024. 1
- 668 Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv*
669 *preprint arXiv:1608.03983*, 2016. 7
670
- 671 David JC MacKay. Information-based objective functions for active data selection. *Neural compu-*
672 *tation*, 4(4):590–604, 1992. 1
- 673 Zhiping Mao and Xuhui Meng. Physics-informed neural networks with residual/gradient-based
674 adaptive sampling methods for solving partial differential equations with sharp solutions. *Applied*
675 *Mathematics and Mechanics*, 44(7):1069–1084, 2023. 6
676
- 677 James D Murray. *Mathematical biology: I. An introduction*, volume 17. Springer Science & Busi-
678 ness Media, 2007. 1
- 679 Daniel Musekamp, Marimuthu Kalimuthu, David Holzmüller, Makoto Takamoto, and Mathias
680 Niepert. Active learning for neural pde solvers. *arXiv preprint arXiv:2408.01536*, 2024. 1,
681 3, 4, 5, 6, 7, 8
682
- 683 Raphaël Pestourie, Youssef Mroueh, Thanh V Nguyen, Payel Das, and Steven G Johnson. Ac-
684 tive learning of deep surrogates for pdes: application to metasurface design. *npj Computational*
685 *Materials*, 6(1):164, 2020. 1, 6
- 686 Ethan Pickering, Stephen Guth, George Em Karniadakis, and Themistoklis P Sapsis. Discovering
687 and forecasting extreme events via active learning in neural operators. *Nature Computational*
688 *Science*, 2(12):823–833, 2022. 1, 6, 8
689
- 690 Harvey M Salkin and Cornelis A De Kluyster. The knapsack problem: a survey. *Naval Research*
691 *Logistics Quarterly*, 22(1):127–144, 1975. 5
- 692 Alvaro Sanchez-Gonzalez, Jonathan Godwin, Tobias Pfaff, Rex Ying, Jure Leskovec, and Peter
693 Battaglia. Learning to simulate complex physics with graph networks. In *International conference*
694 *on machine learning*, pp. 8459–8468. PMLR, 2020. 1
695
- 696 Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set
697 approach. *arXiv preprint arXiv:1708.00489*, 2017. 5, 6
- 698 Burr Settles. Active learning literature survey. Technical report, University of Wisconsin-Madison
699 Department of Computer Sciences, 2009. 1
700
- 701 H Sebastian Seung, Manfred Opper, and Haim Sompolinsky. Query by committee. In *Proceedings*
of the fifth annual workshop on Computational learning theory, pp. 287–294, 1992. 3, 5, 6

- 702 Rohit Singh, Nathan Palmer, David Gifford, Bonnie Berger, and Ziv Bar-Joseph. Active learning for
703 sampling in time-series experiments with application to gene expression analysis. In *Proceedings*
704 *of the 22nd international conference on Machine learning*, pp. 832–839, 2005. 6
- 705
706 Freddie Bickford Smith, Andreas Kirsch, Sebastian Farquhar, Yarin Gal, Adam Foster, and Tom
707 Rainforth. Prediction-oriented bayesian active learning. In *International Conference on Artificial*
708 *Intelligence and Statistics*, pp. 7331–7348. PMLR, 2023. 5
- 709 Karl E Taylor, Ronald J Stouffer, and Gerald A Meehl. An overview of cmip5 and the experiment
710 design. *Bulletin of the American meteorological Society*, 93(4):485–498, 2012. 16
- 711
712 Gabriel Turinici. Optimal time sampling in physics-informed neural networks. *arXiv preprint*
713 *arXiv:2404.18780*, 2024. 6
- 714 Paul Wilmott, Sam Howison, and Jeff Dewynne. *The mathematics of financial derivatives: a student*
715 *introduction*. Cambridge university press, 1995. 1
- 716
717 Chenxi Wu, Min Zhu, Qinyang Tan, Yadhu Kartha, and Lu Lu. A comprehensive study of non-
718 adaptive and residual-based adaptive sampling for physics-informed neural networks. *Computer*
719 *Methods in Applied Mechanics and Engineering*, 403:115671, 2023a. 6
- 720 Dongrui Wu. Pool-based sequential active learning for regression. *IEEE transactions on neural*
721 *networks and learning systems*, 30(5):1348–1359, 2018. 1, 5
- 722
723 Dongxia Wu, Ruijia Niu, Matteo Chinazzi, Yian Ma, and Rose Yu. Disentangled multi-fidelity deep
724 bayesian active learning. In *International Conference on Machine Learning*, pp. 37624–37634.
725 PMLR, 2023b. 1, 6
- 726 Norman J Zabusky and Martin D Kruskal. Interaction of " solitons" in a collisionless plasma and the
727 recurrence of initial states. *Physical review letters*, 15(6):240, 1965. 7, 15
- 728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

756 A EXPERIMENTAL DETAILS

757
758 A.1 DETAILS ON PDES

759 In this section, we describe the PDEs used in our experiments. Each of these equations plays a
760 critical role in modeling physical phenomena and showcases diverse behaviors, from diffusion and
761 soliton dynamics to chaotic systems and fluid flow. Examples of PDE trajectories are shown in
762 Fig. 6.
763



783 **Figure 6:** Example trajectories of PDEs. (a), (b), (c): Horizontal and vertical axes represent the
784 temporal and spatial domain. (d): Two-dimensional states at six time points are shown.

785
786 **Heat Equation** The one-dimensional (1D) Heat equation is given by:

787
788
$$\partial_t u = \partial_{xx} u, \tag{7}$$

789 where $u = u(x, t)$ represents the temperature distribution as a function of space x and time t . This
790 equation describes the process of heat conduction and diffusion in a medium. The simplicity of
791 the Heat equation makes it a fundamental model for understanding diffusion-like processes across
792 various fields in science and engineering, such as thermal conduction, population dynamics, and
793 chemical diffusion. For our experiments, we solve this equation using the pseudospectral method
794 with Dormand–Prince solver as in Brandstetter et al. (2022a).

795
796 **Korteweg–De Vries (KdV) Equation** The second equation we study is the Korteweg–De Vries
797 (KdV) equation, given by:

798
$$\partial_t u + u \partial_x u + \partial_{xxx} u = 0, \tag{8}$$

799 where $u = u(x, t)$ represents a wave profile evolving over space and time. This nonlinear PDE
800 describes the evolution of shallow water waves, and its most famous characteristic is the presence of
801 solitons—solitary, stable wave packets that maintain their shape over long distances and weak inter-
802 actions with other waves (Zabusky & Kruskal, 1965). Solitons have important applications in fluid
803 dynamics, plasma physics, and optical fiber communications. The KdV equation’s nonlinearity and
804 third-order spatial derivative (∂_{xxx}) allow it to capture complex wave behavior. The equation is also
805 known for conserving key quantities like energy. We solve this equation using the pseudospectral
806 method with Dormand–Prince solver as in Brandstetter et al. (2022a).

807
808 **Kuramoto–Sivashinsky (KS) Equation** The Kuramoto–Sivashinsky (KS) equation is a fourth-
809 order nonlinear PDE, written as:

$$\partial_t u + \partial_{xx} u + \partial_{xxxx} u + u \partial_x u = 0, \tag{9}$$

Table 5: Domain lengths and discretizations for trajectory learning.

PDE	Domain Length (T, X)	Resolution (L, N_x)
Heat	(13.0, 6.28)	(13, 256)
KdV	(52.0, 128.0)	(13, 256)
KS	(13.0, 1.0)	(26, 256)
NS	(13.0, 1.0, 1.0)	(13, 32, 32)

where $u = u(x, t)$ is the evolving field in space and time. The KS equation is known for its chaotic behavior and is used to model phenomena such as flame front propagation, plasma instabilities, and thin film dynamics. Its chaotic nature arises from the interplay between destabilizing nonlinear terms and stabilizing higher-order diffusion terms. The equation is particularly challenging to solve due to its sensitivity to initial conditions and long-term unpredictability. To handle this complexity, we use the Exponential Time Differencing (ETD) fourth-order Runge-Kutta method, as introduced by Kassam & Trefethen (2005). This numerical method is well-suited for stiff PDEs like the KS equation.

Navier-Stokes (NS) Equation The final equation we consider is the vorticity form of the incompressible Navier-Stokes (NS) equation, which governs the motion of viscous fluid flows. In two spatial dimensions, the vorticity formulation is given by:

$$\partial_t u + \mathbf{v} \cdot \nabla u = \nu \nabla^2 u + f, \quad \nabla \cdot \mathbf{v} = 0, \quad (10)$$

where $u(x_1, x_2, t)$ is the vorticity, \mathbf{v} is the velocity field, ν is the kinematic viscosity, and $f(x_1, x_2)$ is an external forcing term. The Navier-Stokes equations describe the behavior of incompressible fluid flow, playing a central role in understanding turbulence, weather patterns, and aerodynamics. The external forcing term $f(x_1, x_2)$ is set to

$$f(x) = 0.1 (\sin(2\pi(x_1 + x_2)) + \cos(2\pi(x_1 + x_2))), \quad (11)$$

which injects energy into the system, driving complex fluid dynamics. In our experiments, we adapt the Crank–Nicolson method implemented by Li et al. (2020).

Initial conditions As per Brandstetter et al. (2022a), states are first sampled from a simple distribution and then evolved for a certain time to obtain the initial conditions. The evolved initial conditions are more realistic than the sampled states, in that they are more likely to be observed under a system governed by the respective PDEs. This procedure hence approximates applications where the initial conditions of interest are realistic states either from observed data (Jumper et al., 2021; Kalnay, 2003; Chassignet et al., 2007; Taylor et al., 2012) or carefully crafted synthetic data (Jarrin et al., 2006; Kusner et al., 2017). For 1D equations, the states are sampled from truncated Fourier series with random coefficients (Brandstetter et al., 2022a), and for the 2D NS equation, states are sampled from a Gaussian random field as described in Li et al. (2020). The lengths and discretizations of trajectories are summarized in Table 5.

A.2 ERROR METRICS

The test set always consists of 1,000 trajectories, on which several error metrics are defined. The **RMSE** is defined on a trajectory \mathbf{u} as

$$\sqrt{\frac{1}{LN_x} \sum_{i=1}^L \sum_{j=1}^{N_x} \|\mathbf{u}^i(\mathbf{x}_j) - \hat{\mathbf{u}}^i(\mathbf{x}_j)\|_2^2}. \quad (12)$$

Similarly, the **NRMSE** is defined as

$$\sqrt{\frac{\sum_{i,j} \|\mathbf{u}^i(\mathbf{x}_j) - \hat{\mathbf{u}}^i(\mathbf{x}_j)\|_2^2}{\sum_{i,j} \|\mathbf{u}^i(\mathbf{x}_j)\|_2^2}} \quad (13)$$

Table 6: Acquired datasize in KdV

Round	0	1	2	3	4	5	6	7	8	9	10
SBAL	416	520	624	728	832	936	1040	1144	1248	1352	1456
SBAL+FLEXAL	416	507	611	715	819	923	1027	1131	1235	1339	1443

and the MAE as

$$\frac{1}{LN_x} \sum_{i=1}^L \sum_{j=1}^{N_x} |\mathbf{u}^i(\mathbf{x}_j) - \hat{\mathbf{u}}^i(\mathbf{x}_j)|. \quad (14)$$

The metrics are averaged across all trajectories in the test set. We also report their logarithmic values averaged across all AL rounds, following [Holzmüller et al. \(2023\)](#). Note that we do not use a committee’s mean prediction for computing the metrics, but instead compute the metrics for each model and report their average.

A.3 SIMULATION INSTABILITY

It was observed that using FLEXAL on the KdV equation, the simulation crashes on a small subset of synthetic inputs. Analysis reveals that these synthetic inputs have unusually large norms and particularly appear in later parts of trajectories due to accumulated error. We do not attempt to fix this problem explicitly due to the risk of over-complicating our method, and simply refrain from adding these time steps to the training dataset. This means that FLEXAL actually acquires a smaller number of time steps than the budget B per round of acquisition, which could be problematic when a large subset of inputs do crash. However, we find that this is not the case, and the number of such inputs is small enough that FLEXAL can outperform other baselines. We report the comparison of datasize across rounds in [Table 6](#), for a single experiment. We can see that 13 time steps were left out in the first round due to instability, and no instability occurred in the rounds after.

Since queries that crash incur a cost, they should be avoided as much as possible. Previous works in Bayesian optimization ([Gelbart et al., 2014](#); [Hernández-Lobato et al., 2015](#)) propose methods to learn these unknown constraints. Alternatively, one could simply test out large, random inputs. In fact, we find that the maximum absolute value of an input being above 10 is a robust criterion for predicting that the solver will crash. Either way, we could simply filter out time steps that fall outside of these constraints during runtime of the solver, and use the freed up budget on acquiring other trajectories. Another possible approach is to impose physical constraints on the surrogate model ([Goswami et al., 2022](#)) that reduces the risk of outputting abnormal synthetic inputs. For instance, the KdV equation is energy-conserving, and when this prior knowledge is encoded into the surrogate model, the synthetic inputs would never be abnormally large like we experienced with our naive surrogate models.

A.4 FLEXAL MF

We can also define a simpler acquisition function in the spirit of mean-field approximation. We take the mean model $\hat{G} = \frac{1}{M} \sum_m \hat{G}_m$, and define the variance reduction $R(\hat{G}, b, S)$ between \hat{G} and a model \hat{G}_b in the same way as before. We then average the variance reduction between the mean model and all models in the committee:

$$a_{\text{MF}}(\mathbf{u}^0, S) = \frac{1}{M} \sum_{b \in [M]} R(\hat{G}, b, S), \quad (15)$$

which reduces the computational cost by a factor of M in the best case. We call this modified version FLEXAL MF.

B FULL REPORT OF RESULTS ON MAIN EXPERIMENT

We provide a full report of all results from the main experiment. [Table 7](#), [Table 8](#), [Table 9](#), [Table 10](#) show the full results on Heat, KdV, KS, and NS equations, respectively. [Fig. 7](#) shows the plots of RMSE quantiles on all PDEs. [Fig. 8](#) shows the plots of NRMSE on all PDEs.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

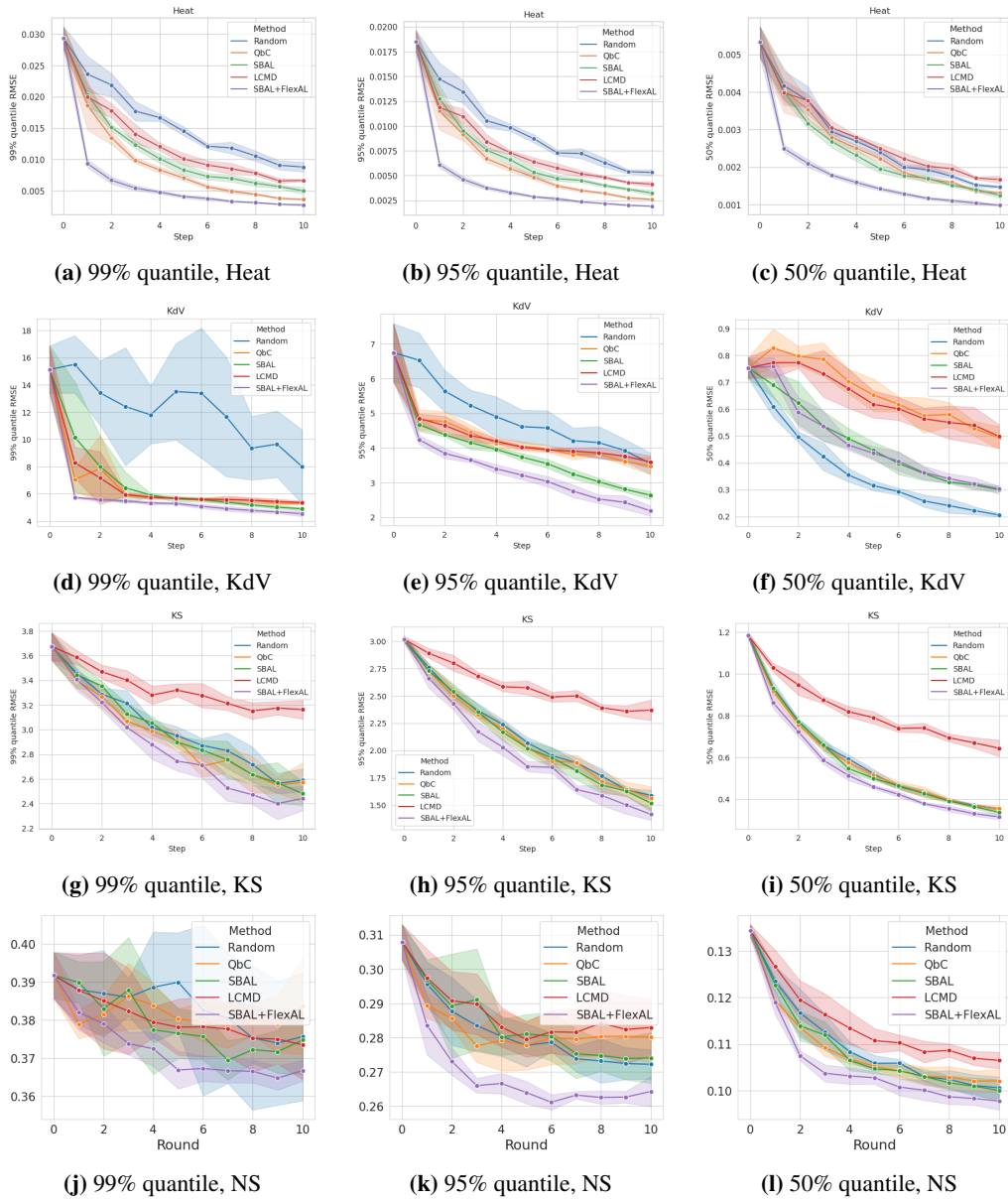


Figure 7: Mean logarithmic values of RMSE quantiles

Table 7: Mean log metrics for Heat Equation

	RMSE	NRMSE	MAE	99%	95%	50%
Random	-5.688 ± 0.021	-6.486 ± 0.018	-7.486 ± 0.021	-4.211 ± 0.029	-4.712 ± 0.022	-5.992 ± 0.020
SBAL	-5.901 ± 0.017	-6.644 ± 0.015	-7.699 ± 0.018	-4.624 ± 0.049	-5.073 ± 0.035	-6.111 ± 0.008
LCMD	-5.741 ± 0.024	-6.494 ± 0.023	-7.541 ± 0.024	-4.466 ± 0.027	-4.945 ± 0.026	-5.940 ± 0.023
QbC	-5.924 ± 0.025	-6.637 ± 0.025	-7.724 ± 0.024	-4.848 ± 0.027	-5.222 ± 0.024	-6.068 ± 0.024
SBAL+FLEXAL	-6.304 ± 0.015	-7.014 ± 0.015	-8.114 ± 0.015	-5.284 ± 0.020	-5.653 ± 0.012	-6.433 ± 0.014
Random+FLEXAL	-6.193 ± 0.021	-6.953 ± 0.017	-7.997 ± 0.021	-4.862 ± 0.037	-5.348 ± 0.032	-6.426 ± 0.016
QbC+FLEXAL	-6.195 ± 0.015	-6.880 ± 0.018	-8.008 ± 0.015	-5.401 ± 0.010	-5.654 ± 0.013	-6.273 ± 0.017
LCMD+FLEXAL	-6.132 ± 0.024	-6.843 ± 0.023	-7.944 ± 0.024	-5.147 ± 0.033	-5.512 ± 0.027	-6.247 ± 0.023

Table 8: Mean log metrics for KdV Equation

	RMSE	NRMSE	MAE	99%	95%	50%
Random	0.191 ± 0.058	-1.193 ± 0.050	-2.034 ± 0.045	2.449 ± 0.047	1.395 ± 0.049	-1.196 ± 0.043
SBAL	0.030 ± 0.029	-1.282 ± 0.030	-2.139 ± 0.027	1.875 ± 0.039	1.267 ± 0.028	-1.267 ± 0.027
QbC	0.266 ± 0.027	-1.019 ± 0.029	-1.879 ± 0.029	1.859 ± 0.037	1.251 ± 0.029	-1.019 ± 0.031
LCMD	0.256 ± 0.030	-1.033 ± 0.036	-1.879 ± 0.033	1.868 ± 0.034	1.322 ± 0.034	-1.100 ± 0.038
SBAL+FLEXAL	-0.088 ± 0.040	-1.378 ± 0.040	-2.239 ± 0.043	1.731 ± 0.040	1.280 ± 0.043	-1.378 ± 0.040
Random+FLEXAL	-0.067 ± 0.054	-1.425 ± 0.047	-2.228 ± 0.036	1.885 ± 0.033	1.296 ± 0.038	-1.424 ± 0.044
QbC+FLEXAL	0.134 ± 0.035	-1.130 ± 0.037	-2.004 ± 0.035	1.721 ± 0.031	1.120 ± 0.037	-1.286 ± 0.035
LCMD+FLEXAL	0.286 ± 0.034	-0.978 ± 0.034	-1.824 ± 0.039	1.799 ± 0.036	1.128 ± 0.034	-1.129 ± 0.032

Following [Holzmüller et al. \(2023\)](#), we also report the 99%, 95%, and 50% quantiles of RMSE. This is useful for analyzing the behavior of AL strategies. AL methods tend to improve performance on points with extreme errors, thus improving performance significantly in the top quantiles, while not so much in the middle quantiles. This is why AL methods perform differently depending on the nature of problem. For instance, problems with more irregularities tend to benefit significantly more from AL methods, since the top quantile errors contribute significantly to the average error in those problems.

As expected, the baseline methods improve performance over random sampling in the 99% quantile, but not so much in the 95% and 50% quantiles. Surprisingly, FLEXAL robustly outperforms the baselines in all error quantiles, which is rarely the case for existing AL methods. We can therefore infer that FLEXAL isn't simply sacrificing the surrogate model's performance in some trajectories to improve its performance in others. FLEXAL both sees a more diverse set of trajectories, and samples the most informative time steps in each trajectory, effectively accounting for how it can improve performance in both the high and middle quantiles of error.

C ADDITIONAL EXPERIMENTS

C.1 DIVERSITY OF SPARSELY SELECTED TIME STEPS

We provide a simple analysis to show that time steps sampled in a sparse manner are more diverse than time steps from entire trajectories. Out of 128 trajectories, we first randomly chose 10 trajectories, which contains $L \times 10$ states. Then, out of all $L \times 128$ states, we randomly chose $L \times 10$ states. The first choice represents full trajectory sampling, and the latter represents sparse time steps sampling. We probe an FNO surrogate model trained on all the 128 trajectories at its hidden layer, and observe the hidden layer activation at each of the $L \times 128$ states. The result is shown in [Fig. 9](#), where black points represent states from the fully sampled trajectories and red points represent sparsely selected states. The latter states are visibly more diverse, which partially explains how sampling time steps in a sparse manner from trajectories can benefit a surrogate model.

Table 9: Mean log metrics for KS Equation

	RMSE	NRMSE	MAE	99%	95%	50%
Random	-0.258 ± 0.004	-1.683 ± 0.004	-2.165 ± 0.004	1.097 ± 0.003	0.752 ± 0.005	-0.575 ± 0.004
SBAL	-0.275 ± 0.014	-1.700 ± 0.014	-2.184 ± 0.014	1.086 ± 0.017	0.732 ± 0.023	-0.594 ± 0.012
QbC	-0.268 ± 0.004	-1.693 ± 0.004	-2.178 ± 0.004	1.077 ± 0.008	0.739 ± 0.013	-0.582 ± 0.006
SBAL+FLEXAL	-0.349 ± 0.003	-1.774 ± 0.003	-2.265 ± 0.003	1.042 ± 0.011	0.672 ± 0.012	-0.674 ± 0.008
Random+FLEXAL	-0.335 ± 0.014	-1.759 ± 0.014	-2.248 ± 0.014	1.060 ± 0.015	0.691 ± 0.007	-0.662 ± 0.015
QbC+FLEXAL	-0.331 ± 0.014	-1.756 ± 0.014	-2.246 ± 0.014	1.050 ± 0.013	0.681 ± 0.020	-0.650 ± 0.013
LCMD	0.046 ± 0.015	-1.378 ± 0.015	-1.829 ± 0.015	1.204 ± 0.009	0.954 ± 0.011	-0.203 ± 0.016
LCMD+FLEXAL	-0.138 ± 0.017	-1.561 ± 0.016	-2.033 ± 0.017	1.139 ± 0.006	0.841 ± 0.014	-0.431 ± 0.016

Table 10: Mean log metrics for NS Equation

	RMSE	NRMSE	MAE	99%	95%	50%
SBAL	-2.052 ± 0.009	-4.253 ± 0.009	-4.592 ± 0.008	-0.970 ± 0.011	-1.260 ± 0.019	-2.217 ± 0.010
Random	-2.050 ± 0.011	-4.249 ± 0.011	-4.590 ± 0.010	-0.959 ± 0.029	-1.266 ± 0.011	-2.208 ± 0.016
QbC	-2.057 ± 0.009	-4.258 ± 0.009	-4.597 ± 0.009	-0.962 ± 0.006	-1.261 ± 0.019	-2.217 ± 0.011
LCMD	-2.018 ± 0.017	-4.219 ± 0.017	-4.560 ± 0.016	-0.967 ± 0.017	-1.248 ± 0.022	-2.168 ± 0.019
SBAL+FLEXAL	-2.092 ± 0.003	-4.293 ± 0.003	-4.632 ± 0.003	-0.988 ± 0.005	-1.309 ± 0.004	-2.249 ± 0.011
Random+FLEXAL	-2.080 ± 0.003	-4.280 ± 0.003	-4.621 ± 0.003	-0.980 ± 0.008	-1.303 ± 0.005	-2.235 ± 0.007
QbC+FLEXAL	-2.079 ± 0.008	-4.280 ± 0.008	-4.619 ± 0.008	-0.979 ± 0.003	-1.293 ± 0.017	-2.238 ± 0.005
LCMD+FLEXAL	-2.051 ± 0.007	-4.253 ± 0.007	-4.593 ± 0.007	-0.984 ± 0.011	-1.258 ± 0.009	-2.211 ± 0.013

C.2 REGULARIZATION FOR TRAJECTORY LEARNING

Brandstetter et al. (2022b) identifies a potential problem with training an autoregressive surrogate model with teacher-forcing. The model experiences a distribution shift during inference, because errors accumulate during rollout unlike during training. They propose a simple fix, called the push-forward trick, which supervises the model \hat{G} not with pairs of \mathbf{u}^{i-1} and \mathbf{u}^i , but with pairs of $\hat{G}[\mathbf{u}^{i-2}]$ and \mathbf{u}^i , where \hat{G} is constantly changing throughout training. An even simpler fix that they experiment with is augmenting the inputs with a Gaussian noise.

One might hypothesize that the advantage of FLEXAL comes from its regularizing effect, since the synthetic inputs in the training set are outputs from the surrogate model \hat{G} . We therefore apply the pushforward trick and Gaussian noise augmentation on the best performing baseline method, SBAL. The results in Table 11 shows that the effect of such regularization methods is minimal compared to the effect of FLEXAL. This shows that the advantage of FLEXAL lies not just in its regularizing effect.

C.3 RANDOM BERNOULLI SAMPLING OF TIME STEPS

We provide the whole list of results with Bernoulli sampling described in Section 5.6. Also, we can enforce consecutive initial time steps sampling by bringing all the true entries in S to the beginning. We call this method Initial Bernoulli sampling, or Initial Ber(p). We report the results with SBAL in Table 12 and Table 13. Initial Bernoulli sampling always performs the worst, possibly because they rarely see the time steps at the end.

C.4 EFFICIENT VARIANTS OF FLEXAL

We provide results with two efficient variants of FLEXAL, namely FLEXAL MF and FLEXAL 10. The results are summarized in Table 14. We provide the wall-clock times of selection algorithms in all equations, in Table 15.

We haven't done extensive experiments with different values of T and ϵ . Increasing T improves performance until it plateaus. Increasing ϵ values higher than a certain point deteriorates the performance slightly. On the other hand, decreasing ϵ too much also deteriorates the performance, but can be recovered with a higher value of T , leading to higher computational cost.

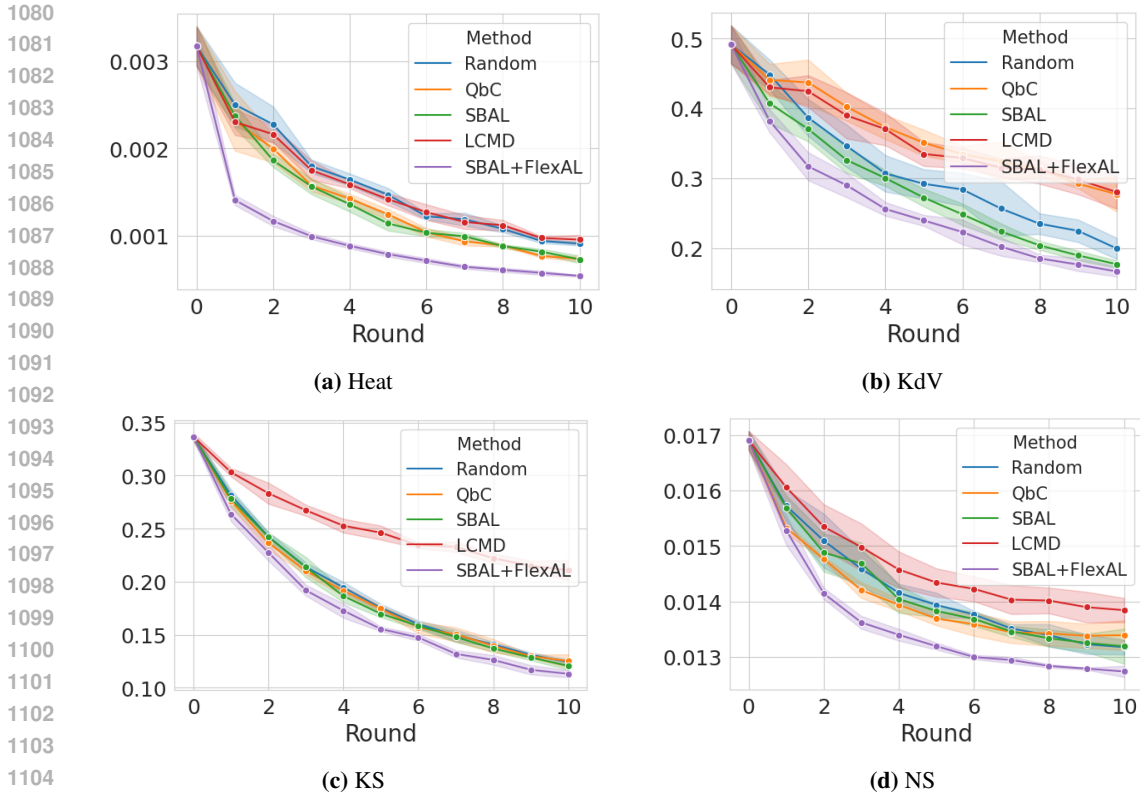


Figure 8: NRMSE of AL strategies, measured across 10 rounds of acquisition. Each round incurs constant cost of data acquisition, namely the budget B . These are simply scaled versions of the RMSE plots.

D FURTHER EXPLANATION OF ACQUISITION WITH FLEXAL

D.1 MOTIVATION BEHIND THE ACQUISITION FUNCTION

Here we detail the motivation behind our acquisition function defined in Section 3.2. First, one can imagine several alternative acquisition functions.

The most straightforward alternative is to simply use the sum of the variances at time points for which $b_i = \text{true}$. The variances are larger for the later time steps since they accumulate, and in our preliminary experiments, we found that this is catastrophic as undersampling the earlier time steps leads to the sampled trajectory being very out-of-distribution, and hence the trained surrogate model underperforming on the test distribution.

It quickly became clear to us that we need some kind of measure of "how much total uncertainty will be reduced by sampling these time steps", instead of "how uncertain is our model on these time steps?" This would help select sampling patterns that reduce the out-of-distribution-ness introduced by \hat{G} . One way to approximate this is to use mutual information, as used by Li et al. (2022b). In other words, we would rollout M trajectories with M surrogate models, and compute the mutual information between time steps for which $b_i = \text{true}$ and all time steps. However, in preliminary experiments, we found that this method underperforms, which we hypothesize is because relying simply on the covariance matrix of the committee between time steps is not a good enough method for computing the posterior uncertainty.

We identified two "pathways" through which sampling a time step reduces uncertainty in the remaining time steps. First, there is the "indirect" pathway: sampling a time step will reduce the model's uncertainty on similar inputs, hence reducing uncertainty on the remaining time steps. This is what is approximated by mutual information. Then, there is the "direct" pathway: sampling a time step

1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187

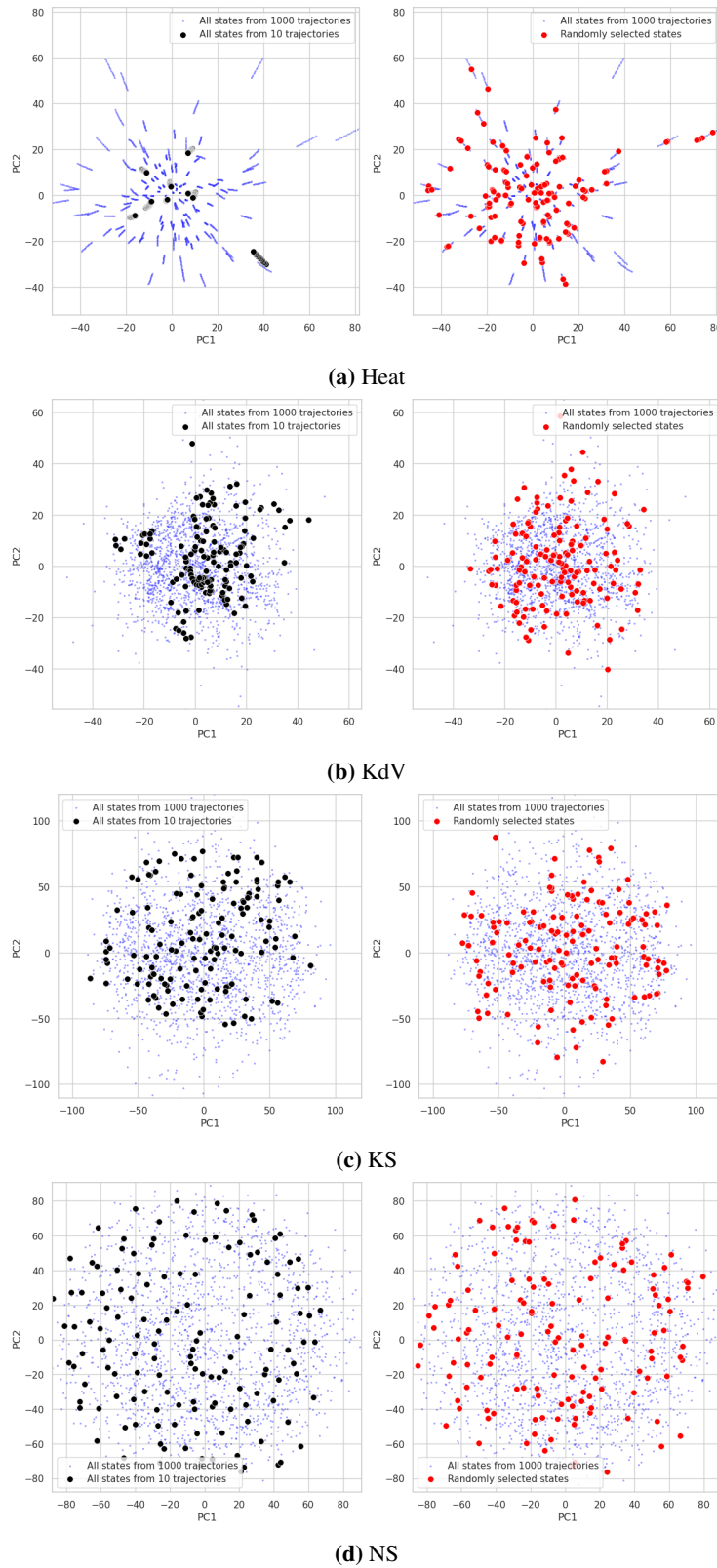


Figure 9: PCA of FNO hidden layer’s activation pattern for both entire trajectories (black) and sparsely sampled time steps (red)

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

Table 11: Effect of regularization

	SBAL	+FLEXAL	+Pushforward	+Gaussian
Heat				
RMSE	-5.901±0.017	-6.304 ±0.015	-3.086±1.584	-5.844±0.011
NRMSE	-6.644±0.015	-7.014 ±0.015	-3.755±1.627	-6.558±0.013
MAE	-7.699±0.018	-8.114 ±0.015	-4.884±1.583	-7.630±0.012
KdV				
RMSE	0.030±0.029	-0.088 ±0.040	0.924±0.613	0.017±0.042
NRMSE	-1.282±0.030	-1.378 ±0.040	-0.245±0.689	-1.292±0.042
MAE	-2.139±0.027	-2.239 ±0.043	-1.077±0.690	-2.162±0.041
KS				
RMSE	-0.275±0.014	-0.349 ±0.003	1.148±0.795	-0.259±0.012
NRMSE	-1.700±0.014	-1.774 ±0.003	-0.283±0.792	-1.684±0.012
MAE	-2.184±0.014	-2.265 ±0.003	-0.473±0.956	-2.167±0.012
NS				
RMSE	-2.052±0.009	-2.092 ±0.003	-0.060±1.118	-2.067±0.012
NRMSE	-4.253±0.009	-4.293 ±0.003	-2.258±1.119	-4.267±0.012
MAE	-4.592±0.008	-4.632 ±0.003	-2.619±1.107	-4.606±0.012

Table 12: Bernoulli sampling

	SBAL	+FLEXAL	+Ber(1/16)	+Ber(1/8)	+Ber(1/4)	+Ber(1/2)
Heat						
RMSE	-5.901±0.017	-6.304 ±0.015	-6.093±0.018	-6.071±0.020	-6.057±0.026	-6.010±0.035
NRMSE	-6.644±0.015	-7.014 ±0.015	-6.823±0.019	-6.801±0.020	-6.791±0.027	-6.748±0.033
MAE	-7.699±0.018	-8.114 ±0.015	-7.893±0.019	-7.872±0.019	-7.858±0.027	-7.810±0.034
KdV						
RMSE	0.030±0.029	-0.088 ±0.040	0.053±0.014	0.049±0.014	0.018±0.024	-0.064±0.031
NRMSE	-1.282±0.030	-1.378 ±0.040	-1.254±0.017	-1.257±0.014	-1.288±0.020	-1.370±0.033
MAE	-2.139±0.027	-2.239 ±0.043	-2.082±0.016	-2.083±0.018	-2.120±0.025	-2.207±0.034
KS						
RMSE	-0.275±0.014	-0.349±0.003	-0.365 ±0.008	-0.359±0.006	-0.346±0.008	-0.324±0.007
NRMSE	-1.700±0.014	-1.774±0.003	-1.790 ±0.008	-1.784±0.006	-1.771±0.008	-1.749±0.007
MAE	-2.184±0.014	-2.265±0.003	-2.282 ±0.007	-2.276±0.006	-2.262±0.009	-2.237±0.009
NS						
RMSE	-2.052±0.009	-2.092 ±0.003	-2.088±0.005	-2.081±0.008	-2.079±0.007	-2.075±0.009
NRMSE	-4.253±0.009	-4.293 ±0.003	-4.288±0.005	-4.282±0.008	-4.279±0.007	-4.276±0.009
MAE	-4.592±0.008	-4.632 ±0.003	-4.626±0.004	-4.620±0.008	-4.617±0.007	-4.614±0.009

i gives out the $i + 1$ th state, which starts a chain reaction of reducing model uncertainty on all successive states. Note that these two pathways are not distinct from a strictly theoretical view, but are rather two ways of approximating uncertainty reduction.

The direct pathway motivated our acquisition function based on variance reduction. In variance reduction, we calculate the posterior uncertainty by rolling out the trajectories with N surrogate models, but collapse into one surrogate model at time steps for which $b_i = \text{true}$. This effectively computes the reduced uncertainty due to the effect of the direct pathway. With experiments, we confirmed that this acquisition function behaves just like we wanted: it is slightly biased towards sampling the earlier time steps, and it chooses an appropriate frequency of time steps to sample that leads to good performance.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

Table 13: Initial Bernoulli sampling

	Initial Ber(1/16)	Initial Ber(1/8)	Initial Ber(1/4)	Initial Ber(1/2)
Heat				
RMSE	-6.278±0.016	-6.254±0.015	-6.182±0.019	-6.080±0.017
NRMSE	-6.989±0.014	-6.966±0.014	-6.902±0.018	-6.811±0.017
MAE	-8.088±0.016	-8.062±0.014	-7.987±0.019	-7.881±0.017
KdV				
RMSE	0.032±0.016	-0.015±0.014	-0.001±0.018	0.011±0.014
NRMSE	-1.278±0.014	-1.321±0.017	-1.303±0.018	-1.294±0.014
MAE	-2.150±0.016	-2.197±0.014	-2.181±0.018	-2.168±0.014
KS				
RMSE	-0.302±0.009	-0.293±0.008	-0.287±0.005	-0.283±0.009
NRMSE	-1.728±0.009	-1.719±0.008	-1.713±0.005	-1.708±0.009
MAE	-2.216±0.008	-2.206±0.008	-2.199±0.007	-2.194±0.010
NS				
RMSE	-2.045±0.016	-2.044±0.014	-2.051±0.019	-2.058±0.014
NRMSE	-4.246±0.014	-4.244±0.014	-4.251±0.019	-4.258±0.014
MAE	-4.596±0.016	-4.594±0.014	-4.598±0.019	-4.602±0.014

Table 14: Log RMSE of more efficient FLEXAL variants averaged across 10 rounds.

	SBAL	+FLEXAL	+FLEXAL MF	+FLEXAL 10
Heat	-5.901±0.017	-6.304 ±0.015	-6.303±0.009	-6.058±0.020
KdV	0.030±0.029	-0.088 ±0.040	-0.065±0.034	-0.118±0.024
KS	-0.275±0.014	-0.349 ±0.003	-0.326±0.004	-0.316±0.009
NS	-2.052±0.009	-2.092 ±0.003	-2.093 ±0.004	-2.078±0.010

D.2 BATCH ACQUISITION ALGORITHM

Algorithm 2 summarizes the batch selection algorithm of FLEXAL. Starting with an empty batch \mathcal{B} , the algorithm repeatedly selects initial conditions and their sampling patterns until reaching the budget limit. It first uses the base active learning method \mathcal{A} to choose an initial condition u^0 . Then, it optimizes which time steps to sample through a greedy procedure: starting with a pattern S that samples all time steps (all true values), it performs T iterations of random mutations. In each iteration, it generates a candidate pattern S' by randomly flipping entries in S with probability ϵ (using a binary mask C where each entry is drawn from a Bernoulli distribution and the XOR operation \oplus). If this new pattern achieves a better value according to the cost-weighted acquisition function a^* , it becomes the current pattern. To ensure the budget isn't exceeded, if adding the current pattern would go over budget, the algorithm truncates it by keeping only enough true values to exactly meet the budget. The pair of initial condition and its optimized sampling pattern (u^0, S) is then added to the batch \mathcal{B} .

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

Equation	QbC	+FlexAL	+FlexAL 10
Heat	10.3	43.0	4.3
KdV	10.6	40.1	4.5
KS	18.1	78.4	8.6
NS	45.5	92.2	10.5

Table 15: Wall-clock times of selection algorithms for all equations

Algorithm 2 Batch Acquisition Algorithm

Require: Budget B , base active learning algorithm \mathcal{A} , probability ϵ , number of iterations T for greedy optimization, pool P of initial conditions, cost function $\text{cost}(\cdot)$ for batches.

Ensure: A batch \mathcal{B} of initial conditions and sampling patterns.

```

1:  $\mathcal{B} \leftarrow \emptyset$ 
2: while  $\text{cost}(\mathcal{B}) < B$  do
3:   Acquire an initial condition  $\mathbf{u}^0$  with  $\mathcal{A}$ .
4:   Initialize  $S \leftarrow (\text{true}, \dots, \text{true})$ .
5:   for  $i = 1$  to  $T$  do
6:      $C = (C_1, \dots, C_L)$  where  $C_1, \dots, C_L \stackrel{\text{i.i.d.}}{\sim} \text{Ber}(\epsilon)$ .
7:      $S' = S \oplus C$ 
8:     if  $a^*(\mathbf{u}^0, S') \geq a^*(\mathbf{u}^0, S)$  then
9:        $S \leftarrow S'$ .
10:    end if
11:  end for
12:  if  $\|\mathcal{B}\| + \text{cost}(\mathcal{B}) > B$  then
13:    Keep only the first  $(B - \text{cost}(\mathcal{B}))$  trues from  $S$  and flip the remaining trues.
14:  end if
15:   $\mathcal{B} \leftarrow \mathcal{B} \cup \{(\mathbf{u}^0, S)\}$ .
16: end while

```
