
Deep Multi-Fidelity Active Learning of High-Dimensional Outputs

Shibo Li, Zheng Wang, Robert M. Kirby, Shandian Zhe

School of Computing, University of Utah

Salt Lake City, UT 84112

{shibo, wzhut, kirby, zhe}@cs.utah.edu

Abstract

Many applications, such as in physical simulation and engineering design, demand we estimate functions with high-dimensional outputs. To reduce the expensive cost of generating training examples, we usually choose several fidelities to enable a cost/quality trade-off. In this paper, we consider the active learning task to automatically identify the fidelities and training inputs to query new examples so as to achieve the best learning benefit-cost ratio. To this end, we propose DMFAL, a Deep Multi-Fidelity Active Learning approach. We first develop a deep neural network based multi-fidelity model for high-dimensional outputs, which can flexibly capture strong complex correlations across the outputs and fidelities to enhance the learning of the target function. We then propose a mutual information based acquisition function that extends the predictive entropy principle. To overcome the computational challenges caused by large output dimensions, we use the multi-variate delta method and moment-matching to estimate the output posterior, and Weinstein-Aronszajn identity to calculate and optimize the acquisition function. We show the advantage of our method in several applications of computational physics and engineering design. The code is available at <https://github.com/shib01i/DMFAL>.

1 INTRODUCTION

Many applications require us to compute a mapping from low-dimensional inputs to high-dimensional

outputs. For example, topology optimization (Rozvany, 2009) aims to find an optimal structure (high-dimensional output) given several design parameters (low-dimensional input). Physical simulation uses numerical solvers to solve partial differential equations (PDEs), which maps the PDE parameters and parameterized initial/boundary conditions (low dimensional input) to the high-dimensional solution field on a mesh. The exact computation of these mappings is often costly. Hence, learning a surrogate model to directly predict the output (rather than computing from scratch every time) is of great interest and importance (Kennedy and O’Hagan, 2000; Conti and O’Hagan, 2010).

However, collecting training examples is a bottleneck, because to obtain these examples we still have to conduct the original, expensive computation (*e.g.*, running numerical solvers). To reduce the cost, we can compute the training examples at different fidelities to enable a trade-off between the cost and quality. Low fidelity examples are cheap to acquire but inaccurate while high-fidelity examples are much more accurate yet expensive. Despite the disparity in accuracy, the outputs at different fidelities are strongly correlated. Hence, examples from different fidelities can all be useful in learning the target mapping.

To reduce the cost while maximizing the learning performance, we develop DMFAL, a deep multi-fidelity active learning approach that can identify both the fidelity and input location to query (or generate) new training examples so as to achieve the best benefit-cost ratio. To our knowledge, this is the first work that incorporates multi-fidelity queries in active learning of high-dimensional outputs (See Fig. 3 in the Appendix for more illustrations). Our work naturally extends the standard active learning, where the query cost is assumed to be uniform, and minimizing the total cost is equivalent to minimizing the number of queries.

Specifically, we first propose an expressive deep multi-fidelity model. We use a chain of neural networks (NNs) to model the outputs in each fidelity. Each NN generates a low-dimensional latent output first, and

Proceedings of the 25th International Conference on Artificial Intelligence and Statistics (AISTATS) 2022, Valencia, Spain. PMLR: Volume 151. Copyright 2022 by the author(s).

then projects it to the high-dimensional observational space. Both the original input and latent output are fed into the NN of the next fidelity so that we can efficiently propagate information throughout the fidelities and flexibly capture their complex relationships. Second, we propose an acquisition function based on the mutual information of the outputs between each fidelity and the highest fidelity (at which we predict the target function). This can be viewed as an extension of the predictive uncertainty principle for the traditional, single-fidelity active learning. When we seek to query with the highest fidelity, the acquisition function is reduced to the output entropy. We found empirically our acquisition function outperforms the popular BALD (Houlsby et al., 2011) and predictive variance principle (Gal et al., 2017) adapted to the multi-fidelity setting. Third, we address the challenges of computing and optimizing the acquisition function. Due to the large output dimension, it is very expensive or even infeasible to estimate the required covariance and cross covariance matrices with popular Monte-Carlo (MC) Dropout (Gal and Ghahramani, 2016) samples. To overcome this problem, we consider the NN weights in the latent output layers as random variables and all the other weights as hyper-parameters. We develop a stochastic structural variational learning algorithm to jointly estimate the hyper-parameters and posterior of the random weights. We then use the multi-variate delta method to compute the moments of the hidden outputs in each fidelity, and use moment-matching to estimate a joint Gaussian posterior of the hidden outputs. We use Weinstein-Aronszajn identity to compute the entropy of their projection — the observed high-dimensional outputs. In this way, we can analytically calculate and optimize the acquisition function in a tractable, reliable and efficient way.

For evaluation, we examined DMFAL in three benchmark tasks of computational physics, a topology structure optimization problem, and a computational dynamic fluids (CFD) application to predict flow velocity fields. The output dimensions of these applications vary from hundreds to hundreds of thousands. Our method consistently achieves much better learning performance with the same query cost, as compared with using random query strategies, approximating the acquisition function with dropout samples of the latent outputs, and using other acquisition functions. Our learned surrogates gain 40x and 460x speed-up in solving (predicting) optimal structures and velocity fields than running standard numerical methods. Even counting the total cost of active training, the per-solution cost of DMFAL rapidly becomes far below that of the numerical methods with the increase of acquired solutions.

2 BACKGROUND

Problem setting. We assume that we can query training examples with M fidelities, which correspond to M mappings, $\{\mathbf{y}_m(\mathbf{x}) \in \mathbb{R}^{d_m}\}_{1 \leq m \leq M}$ where $\mathbf{x} \in \mathcal{X}$ is an r -dimensional input, r is small, d_m is the output dimension, often very large, and $\{d_m\}$ are not necessarily identical. In general, we assume $d_1 \leq \dots \leq d_M$. For instance, in structure design, the input can be the design parameters. The high fidelity examples correspond to high-resolution structures while the low fidelity ones low-resolution structures. We denote by λ_m the cost to query with fidelity m . Our goal is to estimate $\mathbf{y}_M(\mathbf{x})$ (*i.e.*, the target function). Among the examples at different fidelities are strong (but complex) correlations, due to the underlying common bases, say, in physics. Hence, all these examples are valuable to estimate the target function. To perform active learning, we begin with a small set of examples with mixed fidelities. Each step, we select both the fidelity and input location to query (or generate) a new example, so as to best balance the learning improvement (on $\mathbf{y}_M(\mathbf{x})$) and computational cost, *i.e.*, maximizing the benefit-cost ratio.

Dropout active learning. A successful application of active learning with deep neural networks is image classification (Gal et al., 2017). Typically, a pool of unlabeled input examples, *i.e.*, images, is pre-collected. Each time, we rank the input examples according to an acquisition function computed based on the current model. We then query the labels of the top examples and add them into the training set. There are two popular acquisition functions. The first one is the predictive entropy of the label,

$$\mathbb{H}[y|\mathbf{x}, \mathcal{D}] = - \sum_c p(y = c|\mathbf{x}, \mathcal{D}) \log p(y = c|\mathbf{x}, \mathcal{D}), \quad (1)$$

where \mathbf{x} is the input, and \mathcal{D} are the training data. The other one is BALD (Bayesian active learning by disagreement) (Houlsby et al., 2011), namely the mutual information between the label and model parameters,

$$\mathbb{I}[y, \boldsymbol{\theta}|\mathbf{x}, \mathcal{D}] = \mathbb{H}[y|\mathbf{x}, \mathcal{D}] - \mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D})} [\mathbb{H}(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}, \mathcal{D})], \quad (2)$$

where $\boldsymbol{\theta}$ are the NN weights. To calculate these acquisition functions, we need to first estimate the posterior of the NN output and weights. A popular approach is to use Monte-Carlo (MC) Dropout (Gal and Ghahramani, 2016), which is essentially a variational inference method of Bayesian NNs. MC Dropout uses a probability p_i to randomly drop the neurons in each layer i (equivalent to zeroing out the corresponding weights) in the forward pass and performs backward pass over the remaining neurons. After training, the weights and output obtained from one such forward pass can be viewed as a sample of the approximate posterior. We

can run Dropout multiple times to collect a set of posterior samples and calculate the acquisition function by Monte-Carlo approximations.

3 DEEP MULTI-FIDELITY MODELING FOR HIGH-DIMENSIONAL OUTPUTS

Despite its success, dropout active learning might be inappropriate for high-dimensional (continuous) outputs. When the output is a continuous vector, it is natural to use a multivariate Gaussian posterior, $p(\mathbf{y}|\mathbf{x}, \mathcal{D}) \approx \mathcal{N}(\mathbf{y}|\boldsymbol{\mu}(\mathbf{x}), \boldsymbol{\Sigma}(\mathbf{x}))^1$, and the entropy term $\mathbb{H}(\mathbf{y}|\mathbf{x}, \mathcal{D})$ in the acquisition functions (see (1) and (2)) require us to compute the log determinant of the covariance matrix, $\log|\boldsymbol{\Sigma}(\mathbf{x})|$. While we can use dropout samples to construct an empirical estimate $\widehat{\boldsymbol{\Sigma}}(\mathbf{x})$, due to the large output dimension n , the computation of this $n \times n$ matrix and its log determinant is very expensive or even infeasible. One might seek to only estimate the predictive variance of each individual output. However, this will ignore the strong output correlations, which is critical to effectively evaluate the uncertainty and calculate the acquisition function.

To overcome these problems for multi-fidelity active learning, we first propose an expressive deep multi-fidelity model, for which we develop a structural variational inference method to capture the posterior dependency of the outputs. We propose a novel acquisition function to allow multi-fidelity queries, and develop an efficient and reliable approach to calculate the acquisition function.

Specifically, we use deep neural networks to build a multi-fidelity high-dimensional output model (see the graphical representation in Fig. 4 of the Appendix) that can flexibly, efficiently capture complex relationships between the outputs and between the fidelities, taking advantage of these relationships to enhance the learning at the highest fidelity (*i.e.*, target function). For each fidelity m , we introduce a neural network (NN) that first generates a k_m dimensional latent output, and then projects it to the d_m dimensional observational space, where $k_m \ll d_m$. The NN is parameterized by $\{\mathbf{W}_m, \boldsymbol{\theta}_m, \mathbf{A}_m\}$ where \mathbf{W}_m is a $k_m \times l_m$ weight matrix in the latent output layer, \mathbf{A}_m a $d_m \times k_m$ projection matrix, and $\boldsymbol{\theta}_m$ the weights in all the other layers. Denote by $\boldsymbol{\xi}_m$ the NN input, by $\mathbf{h}_m(\mathbf{x})$ the latent NN output, and by $\mathbf{y}_m(\mathbf{x})$ the observed output. The model is defined by

$$\begin{aligned} \boldsymbol{\xi}_m &= [\mathbf{x}; \mathbf{h}_{m-1}(\mathbf{x})], \quad \mathbf{h}_m(\mathbf{x}) = \mathbf{W}_m \boldsymbol{\rho}_{\boldsymbol{\theta}_m}(\boldsymbol{\xi}_m), \\ \mathbf{y}_m(\mathbf{x}) &= \mathbf{A}_m \mathbf{h}_m(\mathbf{x}) + \boldsymbol{\epsilon}_m, \end{aligned} \quad (3)$$

¹one can also use a multivariate student t distribution, but the problem remains.

where $\boldsymbol{\xi}_1 = \mathbf{x}$, $\boldsymbol{\rho}_{\boldsymbol{\theta}_m}$ is the second last layer of the NN at fidelity m , of dimension l_m , and $\boldsymbol{\epsilon}_m \sim \mathcal{N}(\boldsymbol{\epsilon}_m|\mathbf{0}, \sigma_m^2 \mathbf{I})$ is an isotropic Gaussian noise. Note that $\boldsymbol{\rho}_{\boldsymbol{\theta}_m}$ can be considered as l_m nonlinear basis functions parameterized by $\boldsymbol{\theta}_m$. Through their combinations via \mathbf{W}_m and \mathbf{A}_m , we can flexibly capture the complex relationships between the elements of \mathbf{y}_m to improve the prediction. Furthermore, the input $\boldsymbol{\xi}_m$ is the concatenation of the original input \mathbf{x} and the latent output \mathbf{h}_{m-1} . The latter can be viewed as a compact (or low-rank) summary of all the information up to fidelity $m-1$. Through an NN mapping, we obtain the latent output $\mathbf{h}_m = \text{NN}(\mathbf{x}, \mathbf{h}_{m-1}(\mathbf{x}))$ — the summary up to fidelity m — and then generate the high dimensional output \mathbf{y}_m . Thereby, we are able to efficiently integrate the information from lower fidelities and grasp the strong, complex relationship between the current and prior fidelities.

We place a standard Gaussian prior over the elements in each \mathbf{W}_m . Similar to (Snoek et al., 2015), we consider all the remaining parameters as hyper-parameters to ease the posterior inference and uncertainty reasoning. Given the training dataset $\mathcal{D} = \{(\mathbf{x}_{nm}, \mathbf{y}_{nm})\}_{n=1}^{N_m}\}_{m=1}^M$, the joint probability of our model is given by

$$\begin{aligned} p(\mathcal{W}, \mathcal{Y}|\mathcal{X}, \Theta, \mathbf{s}) &= \prod_{m=1}^M \mathcal{N}(\text{vec}(\mathbf{W}_m)|\mathbf{0}, \mathbf{I}) \\ &\cdot \prod_{n=1}^{N_m} \mathcal{N}(\mathbf{y}_{nm}|\mathbf{A}_m \mathbf{h}_m(\mathbf{x}_{nm}), \sigma_m^2 \mathbf{I}), \end{aligned} \quad (4)$$

where $\mathcal{W} = \{\mathbf{W}_m\}_{m=1}^M$, $\Theta = \{\boldsymbol{\theta}_m, \mathbf{A}_m\}_{m=1}^M$, $\mathbf{s} = [\sigma_1^2, \dots, \sigma_M^2]$, and $\{\mathcal{X}, \mathcal{Y}\}$ are the inputs and outputs in \mathcal{D} . To estimate the posterior of our model (which is used to calculate the acquisition function and query new examples), we develop a stochastic structural variational learning algorithm. Specifically, for each \mathbf{W}_m , we introduce a multivariate Gaussian variational posterior, $q(\mathbf{W}_m) = \mathcal{N}(\text{vec}(\mathbf{W}_m)|\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$. To ensure the positive definiteness, we parameterize $\boldsymbol{\Sigma}_m$ by its Cholesky decomposition, $\boldsymbol{\Sigma}_m = \mathbf{L}_m \mathbf{L}_m^\top$, where \mathbf{L}_m is a lower triangular matrix. We then assume the posterior $q(\mathcal{W}) = \prod_{m=1}^M q(\mathbf{W}_m)$, and construct a variational model evidence lower bound (ELBO) (Wainwright et al., 2008), $\mathcal{L}(q(\mathcal{W}), \Theta, \mathbf{s}) = \mathbb{E}_q[\log p(\mathcal{W}, \mathcal{Y}|\mathcal{X}, \Theta, \mathbf{s})/q(\mathcal{W})] = -\text{KL}(q(\mathcal{W})\|p(\mathcal{W})) + \sum_{m=1}^M \sum_{n=1}^{N_m} \mathbb{E}_q[\log \mathcal{N}(\mathbf{y}_{nm}|\mathbf{A}_m \mathbf{h}_m(\mathbf{x}_{nm}), \sigma_m^2 \mathbf{I})]$, where $\text{KL}(\cdot\|\cdot)$ is Kullback Leibler divergence and $p(\mathcal{W})$ the prior of \mathcal{W} . We maximize \mathcal{L} to jointly estimate $q(\mathcal{W})$ and the hyper-parameters. While \mathcal{L} is analytically intractable, we use the reparameterization trick (Kingma and Welling, 2013) to generate parameterized samples for each $\text{vec}(\mathbf{W}_m)$: $\boldsymbol{\mu}_m + \mathbf{L}_m \boldsymbol{\iota}_m$ where $\boldsymbol{\iota}_m \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ to obtain a stochastic estimate of \mathcal{L} , and conduct stochastic optimization.

4 MULTI-FIDELITY ACTIVE LEARNING

4.1 Mutual Information based Acquisition Function

We now consider how to perform active learning with multi-fidelity queries. We assume that at each fidelity m , the most valuable training example is the one that can best improve our prediction of the target function, namely, the prediction at the highest fidelity M . To this end, we propose our acquisition function based on the mutual information between the outputs at fidelity m and M ,

$$\begin{aligned} a(\mathbf{x}, m) &= \frac{1}{\lambda_m} \mathbb{I}(\mathbf{y}_m(\mathbf{x}), \mathbf{y}_M(\mathbf{x}) | \mathcal{D}) \\ &= \frac{1}{\lambda_m} (\mathbb{H}(\mathbf{y}_m | \mathcal{D}) + \mathbb{H}(\mathbf{y}_M | \mathcal{D}) - \mathbb{H}(\mathbf{y}_m, \mathbf{y}_M | \mathcal{D})), \end{aligned} \quad (5)$$

where $\lambda_m > 0$ is the cost of querying a training example with fidelity m . When $m = M$, we have $a(\mathbf{x}, M) = \frac{1}{\lambda_M} \mathbb{H}(\mathbf{y}_M(\mathbf{x}) | \mathcal{D})$ — the output entropy. Therefore, our acquisition function is an extension of the popular predictive entropy principle in conventional active learning. At each step, we maximize our acquisition function to identify a pair of fidelity and input location that give the biggest benefit-cost ratio.

4.2 Efficient Acquisition Function Calculation

To maximize the acquisition function (5), a critical challenge is to compute the posterior of the outputs $\{\mathbf{y}_m\}$ in every fidelity, based on our model estimation results, *i.e.*, $p(\mathcal{W} | \mathcal{D}) \approx q(\mathcal{W})$. Due to the high dimensionality of each \mathbf{y}_m and the nonlinear coupling of the latent outputs across the fidelities (see (3)), the computation is challenging and analytically intractable. To address this issue, we consider approximating the posterior of the low dimensional latent output \mathbf{h}_m first, which can be viewed as a function of the random NN weights $\boldsymbol{\Omega}_m = \{\mathbf{W}_1, \dots, \mathbf{W}_m\}$ (given the input \mathbf{x}). We then use multivariate delta method (Oehlert, 1992; Bickel and Doksum, 2015) to compute the moments of \mathbf{h}_m . Specifically, we approximate \mathbf{h}_m with a first-order Taylor expansion,

$$\mathbf{h}_m(\boldsymbol{\Omega}_m) \approx \mathbf{h}_m(\mathbb{E}[\boldsymbol{\Omega}_m]) + \mathbf{J}_m (\boldsymbol{\eta}_m - \mathbb{E}[\boldsymbol{\eta}_m]), \quad (6)$$

where the expectation is under $q(\cdot)$, $\boldsymbol{\eta}_m = \text{vec}(\boldsymbol{\Omega}_m)$, $\mathbf{J}_m = \frac{\partial \mathbf{h}_m}{\partial \boldsymbol{\eta}_m} |_{\boldsymbol{\eta}_m = \mathbb{E}[\boldsymbol{\eta}_m]}$ is the Jacobian matrix at the mean. The rationale is as follows. First, \mathbf{h}_m is linear to \mathbf{W}_m and the second-order derivative is simply $\mathbf{0}$. Second, as the NN output, \mathbf{h}_m is highly nonlinear to the random weights in the previous layers, \mathbf{W}_j ($j < m$). Hence, we can assume the change rate of \mathbf{h}_m (*i.e.*, gradient or Jacobian) has much greater scales than the

posterior covariance of \mathbf{W}_j in the second-order term of the Taylor expansion. Note that $q(\mathbf{W}_j)$ is more informative and hence much more concentrated than the prior $\mathcal{N}(\text{vec}(\mathbf{W}_j) | \mathbf{0}, \mathbf{I})$. The scale of the posterior covariance should be much smaller than 1. Considering both cases, the first-order term can dominate the Taylor expansion and hence we ignore the higher order terms. Our ablation study has confirmed our analysis; see Sec. D of Appendix for details. Based on (6), we can easily calculate the first and second moments of \mathbf{h}_m ,

$$\boldsymbol{\alpha}_m = \mathbb{E}[\mathbf{h}_m] \approx \mathbf{h}_m(\mathbb{E}[\boldsymbol{\Omega}_m]), \quad (7)$$

$$\mathbf{V}_m = \text{cov}[\mathbf{h}_m] \approx \mathbf{J}_m \text{cov}(\boldsymbol{\eta}_m) \mathbf{J}_m^\top, \quad (8)$$

where $\text{cov}(\boldsymbol{\eta}_m) = \text{diag}(\{\text{cov}(\text{vec}(\mathbf{W}_j))\}_{1 \leq j \leq m})$. Then we use moment matching to estimate a joint Gaussian posterior, $q(\mathbf{h}_m) = \mathcal{N}(\mathbf{h}_m | \boldsymbol{\alpha}_m, \mathbf{V}_m)$. Next, according to (3), we can obtain the posterior of the output \mathbf{y}_m ,

$$q(\mathbf{y}_m) = \mathcal{N}(\mathbf{y}_m | \mathbf{A}_m \boldsymbol{\alpha}_m, \mathbf{A}_m \mathbf{V}_m \mathbf{A}_m^\top + \sigma_m^2 \mathbf{I}). \quad (9)$$

The output entropy is $\mathbb{H}(\mathbf{y}_m | \mathcal{D}) = \frac{1}{2} \log |\mathbf{A}_m \mathbf{V}_m \mathbf{A}_m^\top + \sigma_m^2 \mathbf{I}| + d_m \log \sqrt{2\pi e}$. However, directly computing the log determinant of a $d_m \times d_m$ matrix is very expensive. To bypass this challenge, we use the Weinstein-Aronszajn identity (Kato, 2013) to derive

$$\begin{aligned} \mathbb{H}(\mathbf{y}_m | \mathcal{D}) &= \frac{1}{2} \log |\sigma_m^{-2} \mathbf{A}_m \mathbf{V}_m \mathbf{A}_m^\top + \mathbf{I}| + \log (2\pi e \sigma_m^2)^{\frac{d_m}{2}} \\ &= \frac{1}{2} \log |\sigma_m^{-2} \mathbf{A}_m^\top \mathbf{A}_m \mathbf{V}_m + \mathbf{I}| + d_m \log \sqrt{2\pi e \sigma_m^2}. \end{aligned} \quad (10)$$

Now, the log determinant is calculated on a much smaller, $k_m \times k_m$ matrix, which is very cheap and efficient.

Using a similar approach, we can calculate the joint posterior of $\hat{\mathbf{h}}_m = [\mathbf{h}_m; \mathbf{h}_M]$ and then $\hat{\mathbf{y}}_m = [\mathbf{y}_m; \mathbf{y}_M]$. First, we view $\hat{\mathbf{h}}_m$ as a function of all the random weights, $\mathcal{W} = \{\mathbf{W}_1, \dots, \mathbf{W}_M\}$. We use the multivariate delta method and moment matching to estimate a joint Gaussian posterior, $q(\hat{\mathbf{h}}_m) = \mathcal{N}(\hat{\mathbf{h}}_m | \hat{\boldsymbol{\alpha}}_m, \hat{\mathbf{V}}_m)$ where $\hat{\boldsymbol{\alpha}}_m = \hat{\mathbf{h}}_m(\mathbb{E}[\mathcal{W}])$, $\hat{\mathbf{V}}_m = \hat{\mathbf{J}}_m \text{cov}(\boldsymbol{\eta}) \hat{\mathbf{J}}_m^\top$, $\boldsymbol{\eta} = \text{vec}(\mathcal{W})$, $\hat{\mathbf{J}}_m = \frac{\partial \hat{\mathbf{h}}_m}{\partial \boldsymbol{\eta}} |_{\boldsymbol{\eta} = \mathbb{E}[\boldsymbol{\eta}]}$, and $\text{cov}(\boldsymbol{\eta}) = \text{diag}(\{\text{cov}(\text{vec}(\mathbf{W}_j))\}_{1 \leq j \leq M})$. According to our model (3), we can represent

$$\hat{\mathbf{y}}_m = \hat{\mathbf{A}}_m \hat{\mathbf{h}}_m + \hat{\boldsymbol{\epsilon}}_m,$$

where $\hat{\mathbf{A}}_m = \text{diag}(\mathbf{A}_m, \mathbf{A}_M)$ and $\hat{\boldsymbol{\epsilon}}_m = [\boldsymbol{\epsilon}_m; \boldsymbol{\epsilon}_M]$. Therefore, the joint posterior of $\hat{\mathbf{y}}_m$ is

$$q(\hat{\mathbf{y}}_m) = \mathcal{N}(\hat{\mathbf{y}}_m | \hat{\mathbf{A}}_m \hat{\boldsymbol{\alpha}}_m, \hat{\mathbf{A}}_m \hat{\mathbf{V}}_m \hat{\mathbf{A}}_m^\top + \mathbf{S}_m), \quad (11)$$

where $\mathbf{S}_m = \text{diag}(\sigma_m^2 \mathbf{I}_{d_m}, \sigma_M^2 \mathbf{I}_{d_M})$, \mathbf{I}_{d_m} and \mathbf{I}_{d_M} are identity matrices of $d_m \times d_m$ and $d_M \times d_M$, respectively. Again, we use Weinstein-Aronszajn identity to simplify the entropy computation of $\hat{\mathbf{y}}_m$ (*i.e.*, \mathbf{y}_m and \mathbf{y}_M),

$$\begin{aligned} \mathbb{H}(\mathbf{y}_m, \mathbf{y}_M | \mathcal{D}) &= \frac{1}{2} \log |\mathbf{S}_m^{-1} \hat{\mathbf{A}}_m \hat{\mathbf{V}}_m \hat{\mathbf{A}}_m^\top + \mathbf{I}| + \beta_m \\ &= \frac{1}{2} \log |\hat{\mathbf{A}}_m^\top \mathbf{S}_m^{-1} \hat{\mathbf{A}}_m \hat{\mathbf{V}}_m + \mathbf{I}| + \beta_m, \end{aligned} \quad (12)$$

where $\beta_m = d_m \log \sqrt{2\pi e \sigma_m^2} + d_M \log \sqrt{2\pi e \sigma_M^2}$ and the log determinant is computed from a $(k_m + k_M) \times (k_m + k_M)$ matrix, which is cheap and efficient. Note that we can use block matrices to compute $\hat{\mathbf{A}}_m^\top \mathbf{S}_m^{-1} \hat{\mathbf{A}}_m = \text{diag}(\sigma_m^{-2} \mathbf{A}_m^\top \mathbf{A}_m, \sigma_M^{-2} \mathbf{A}_M^\top \mathbf{A}_M)$.

Now, based on (10) and (12), we can calculate our acquisition function (5) in an analytic and deterministic way. For each fidelity m , we maximize $a(\mathbf{x}, m)$ w.r.t to \mathbf{x} to find the optimal input. We can use automatic differential libraries to calculate the gradient and feed it to any optimization algorithm, *e.g.*, L-BFGS. We then use the optimal input at the fidelity that has the largest acquisition function value to query the next training example. Our deep multi-fidelity active learning is summarized in Algorithm 1.

Algorithm 1 DMFAL (\mathcal{D} , T , $\{\lambda_m\}_{m=1}^M$)

- 1: Train the deep multi-fidelity model (4) on the initial dataset \mathcal{D} with stochastic structural variational learning.
- 2: **for** $t = 1, \dots, T$ **do**
- 3: Based on (10) and (12), calculate and optimize the acquisition function (5) to find

$$(\mathbf{x}_t, m_t) = \underset{\mathbf{x} \in \mathcal{X}, 1 \leq m \leq M}{\text{argmax}} a(\mathbf{x}, m).$$

- 4: Query the output \mathbf{y}_t at input \mathbf{x}_t with fidelity m_t .
 - 5: $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{x}_t, \mathbf{y}_t) | m_t\}$.
 - 6: Re-train the deep multi-fidelity model on \mathcal{D} .
 - 7: **end for**
-

4.3 Algorithm Complexity

The time complexity of training our deep multi-fidelity model is $\mathcal{O}(N(\sum_{m=1}^M (k_m l_m)^2 + F))$ where N and F are the total number of training examples and NN parameters, respectively. The space complexity is $\mathcal{O}(F + \sum_{m=1}^M (k_m l_m)^2)$, which is to store the NN parameters, and posterior mean and covariance of each random weight matrix \mathbf{W}_m . The time complexity of calculating the acquisition function is $\mathcal{O}(\sum_{m=1}^M d_m k_m^2 + k_m^3)$. Since $k_m \ll d_m$, the complexity is linear in the output dimensions. Due to the usage of the learned variational posterior $q(\mathcal{W})$, the space complexity of computing the acquisition function is the same as that of training the multi-fidelity model.

5 RELATED WORK

Active learning (AL) is a fundamental machine learning topic (Balcan et al., 2007; Settles, 2009; Balcan et al., 2009; Dasgupta, 2011; Hanneke et al., 2014). Recent research focuses more on deep neural networks. Gal et al. (2017) used Monte-Carlo (MC) Dropout (Gal and Ghahramani, 2016) to perform variational inference for Bayesian neural networks. The dropout samples can be viewed as the posterior samples of the output. These samples are then used to compute an information measurement, such as predictive entropy and BALD (Houlsby et al., 2011), to select unlabeled examples to query. This method has achieved a success in image classification. There are also other works along this line. For example, (Geifman and El-Yaniv, 2017; Sener and Savarese, 2018) selected representative examples based on core-set search. (Gissin and Shalev-Shwartz, 2019) selected maximally indistinguishable samples from the unlabeled pool, and the idea is reminiscent of generative adversarial networks (Goodfellow et al., 2014). (Ducoffe and Precioso, 2018) used adversarial samples to calculate the distance to the decision boundary and select examples accordingly to label. (Ash et al., 2019) measured the uncertainty in terms of the gradient magnitude and selected a disparate batch of inputs in a hallucinated gradient space.

Our work differs from the existing studies in several aspects. First, most methods are pool-based active learning, *i.e.*, a pool of unlabeled examples are collected beforehand, and new examples are only selected from the pool. This is reasonable when the training input is high-dimensional or hard to generate, *e.g.*, image classification. By contrast, our work focuses on learning mappings from low-dimensional inputs to high-dimensional outputs, which are common in physical simulation and engineering design. To find the best training example, we *optimize* the acquisition function in the entire domain rather than limit the search in a discrete set. Second, most active learning methods assume a uniform fidelity (or quality) of the training examples. Our work considers the case that the training examples can be queried with multiple fidelities, resulting in different cost/quality trade-offs. Our goal is to maximize the learning improvement while minimizing the query cost, *i.e.*, maximizing the benefit-cost ratio. This naturally extends the standard active learning, where the total cost is proportional to the number of queries due to the uniform fidelity. The early work of (Settles et al., 2008) empirically studies active learning with nonuniform human labeling costs. To our knowledge, our work is the first multi-fidelity active learning approach for high-dimensional outputs.

Many multi-fidelity (MF) models have been proposed, while their active training approaches are lacking.

These models are often based on Gaussian processes (GPs). For example, Perdikaris et al. (2017) successfully estimated a set of GPs, where each GP predicts the output for one fidelity, where the input comes from the original input and the output of the previous fidelity. Hamelijnck et al. (2019) used Gaussian process regression networks (GPRN) (Wilson et al., 2012), deep GP (Damianou and Lawrence, 2013) and mixture of experts (Rasmussen and Ghahramani, 2002) to develop a multi-task, multi-resolution GP model. Despite their success, these models do not scale to high-dimensional outputs. The most recent work by Wang et al. (2021) overcomes this limitation by proposing a nonlinear coregionalization component and then stacking these components for multi-fidelity modeling. How to conduct active learning for this model, however, remains an open problem — even considering our proposed acquisition function, the computation is challenging due to the complicated couplings of the kernels and weight functions. In our experiments, we have compared the performance of our model with (Wang et al., 2021), and other state-of-the-art high-dimensional GP regressors (that are single-fidelity) in non-active learning. The results confirmed the advantage of our NN based multi-fidelity model (see Sec. C of Appendix).

AL also connects to Bayesian optimization (BO) (Mockus, 2012), while the latter aims to find the function optimum and so the acquisition function definition, computation, and optimizing techniques/challenges are very different. Interestingly, while multi-fidelity AL methods are scarce, MF BO is common, *e.g.*, (Huang et al., 2006; Picheny et al., 2013; Kandasamy et al., 2016; Poloczek et al., 2017; McLeod et al., 2017; Wu and Frazier, 2017). Most MFBO approaches are based on GPs — the classical BO surrogate. The recent MFBO work of Li et al. (2020) uses a chain of NNs to estimate a multi-fidelity surrogate of the black-box objective function, queries new examples based on the maximum-value entropy search (MES) principle (Wang and Jegelka, 2017). From the modeling perspective, its key difference (from our model) is that the NN output at each fidelity is directly fed into the NN for the next fidelity. Therefore, their model will need massive weights to connect multiple NNs for high-dimensional outputs, which can be very inefficient. In addition, the massive outputs of the previous fidelity might dominate the original input (low-dimensional) to the model. Li et al. (2020) developed a recursive one-dimensional quadrature to estimate the posterior of the output at each fidelity. While being suitable for optimizing a single-output black-box function, this technique is difficult to extend to multiple outputs (due to the explosion of the cost for multi-dimensional quadrature), not to mention massive outputs for active learning.

6 EXPERIMENT

6.1 Solving Partial Differential Equations

We first evaluated DMFAL in standard computational physics tasks. Specifically, we used DMFAL to predict the solution fields of three commonly used partial differential equations (PDEs): *Burgers'*, *Poisson's* and *Heat* equations (Olsen-Kettle, 2011). The training examples are collected by running a numerical solver with different meshes. The more the nodes/steps to create the mesh, the higher the fidelity. The input includes the PDE parameters and/or parameterized initial or boundary conditions. The output consists of the solution values at the mesh used in the solver. For example, a 50×50 mesh corresponds to a 2,500 dimensional output vector. For active learning, we considered two-fidelity queries for all the three equations, where the sizes of the corresponding output fields are 16×16 and 32×32 . In addition, we considered a three-fidelity setting for Poisson's equation, denoted by *Poisson-3*, and the sizes of the output fields of the three fidelities are 16×16 , 32×32 , and 64×64 , respectively. For the two-fidelity active learning, we uniformly sampled the inputs, and queried 10 training examples at the first fidelity and 2 at the second fidelity. We used those examples as the initial training dataset. Similarly, for *Poisson-3*, we collected 10, 5 and 2 examples in the first, second, and third fidelity as the initial training set. We generated 500 samples for test, where the test inputs were uniformly sampled from the domain. The outputs are calculated by running the solver with an even denser mesh — 128×128 for Burger's and Poisson's equations, and 100×100 for Heat equation — and interpolating the solution values at the target grid (Zienkiewicz et al., 1977) (this is the standard approach in physical simulation and the accuracy does not change). More details are given in the Appendix. *We ran the solvers at each fidelity for many times and calculated the average running time. We then normalized the average running time to obtain $\lambda_1 = 1$, $\lambda_2 = 3$ and $\lambda_3 = 10$.*

Competing methods. We compared DMFAL with the following active learning approaches. (1) MF-BALD, a straightforward extension of BALD (Houlsby et al., 2011) to integrate multi-fidelity queries. The acquisition function is $a_{\text{MF-BALD}}(\mathbf{x}, m) = \frac{1}{\lambda_m} \mathbb{I}(\mathbf{y}_m(\mathbf{x}), \mathcal{W}|\mathcal{D}) = \frac{1}{\lambda_m} (\mathbb{H}(\mathbf{y}_m(\mathbf{x})|\mathcal{D}) - \mathbb{E}_{p(\mathcal{W}|\mathcal{D})} [\mathbb{H}(\mathbf{y}_m(\mathbf{x})|\mathcal{W}, \mathcal{D})]) = \frac{1}{\lambda_m} (\mathbb{H}(\mathbf{y}_m(\mathbf{x})|\mathcal{D}) - \frac{d_m}{2} \log(2\pi e \sigma_m^2))$. Note that conditioned the NN parameters, the entropy of the observed output \mathbf{y}_m is only determined by the noise variance σ_m^2 . (2) MF-PredVar, a straightforward extension of the popular predictive variance principle, $a_{\text{MF-PredVar}} = \frac{1}{\lambda_m} \frac{1}{d_m} \sum_{j=1}^{d_m} \text{Var}(y_{mj}|\mathcal{D})$. (3)

Dropout-latent, where we use MC dropout (Gal et al., 2017) for variational inference, each time draw 100 dropout samples for the low dimensional latent outputs $\{\mathbf{h}_1(\mathbf{x}), \dots, \mathbf{h}_M(\mathbf{x})\}$ to estimate multi-variate Gaussian posteriors for each \mathbf{h}_m and $\hat{\mathbf{h}}_m = [\mathbf{h}_m; \mathbf{h}_M]$ (via empirical means and covariance matrices), and then follow Section 4.2 to calculate and optimize the acquisition function (5). Note that we have also used MC dropout to outright sample the final, high-dimensional outputs $\{\mathbf{y}_m\}$ and estimate multi-variate Gaussian posteriors to calculate (5) and $a_{\text{MF-BALD}}$. While doing this is much more expensive, the performance did not improve. Instead, it was way worse than Dropout-latent and MF-BALD. See the details in Sec. A.3 of Appendix. (4) MF-Random, where each time we randomly select a fidelity and then an input to query. (5) Random-F1, (6) Random-F2, and (7) Random-F3, where we stick to the first, second and third fidelity, respectively, and randomly sample an input to query each time.

Settings and results. We introduced a two-layer NN for each fidelity and used `tanh` as the activation function². The layer width was chosen from $\{8, 16, 32, 64, 128\}$. We set the same dimension for the latent output in each fidelity and selected it from $\{5, 10, 15, 20\}$. For Dropout-latent, we tuned the dropout rate from $\{0.1, 0.2, 0.3, 0.4, 0.5\}$. All the methods were implemented with PyTorch (Paszke et al., 2019). We used ADAM (Kingma and Ba, 2014) for stochastic optimization, where the learning rate was selected from $\{10^{-4}, 5 \times 10^{-4}, 10^{-3}, 5 \times 10^{-3}, 10^{-2}\}$. We set the number of epochs to 2,000, which is enough for convergence. We conducted five runs for each method, and in each run, we queried 100 examples. We report the average normalized root-mean-square-error (nRMSE) *vs.* the accumulated cost in Fig. 1a-d. The shaded region shows the standard deviation. We can see that at the beginning, all the methods have the same or comparable performance. Along with more queries, DMFAL quickly achieves better prediction accuracy, and continues to outperform all the other methods by a large margin. Therefore, DMFAL can perform much better with the same cost or achieve the same performance with the smallest cost. It is interesting to see that in Fig. 1a, while all the competing approaches have saturated early, DMFAL keeps improving its prediction accuracy.

The inferior performance of MF-BALD and MF-PredVar to DMFAL indicates that the straightforward

²We tried more layers per fidelity, *i.e.*, 3-5, and the overall performance did not improve. In addition, our experience shows that alternative activation functions, like `ReLU` and `LeakyReLU`, are inferior to `tanh`. This is consistent with the choice of typical physics-informed neural networks (Raissi et al., 2019).

extension of BALD and predictive variance can be sub-optimal, which might partly because they miss the influence of the low-fidelity training examples on the final predictions (*i.e.*, at the highest fidelity). How to better generalize BALD and PredVar to the multi-fidelity setting remains an open problem. Overall, these results have demonstrated the advantage of our deep multi-fidelity active learning approach.

6.2 Topology Structure Optimization

Next, we applied DMFAL in topology structure optimization. A topology structure is a layout of materials, *e.g.*, alloy and concrete, in some designated spatial domain. Given the input from the outside environment, *e.g.*, external force, we want to find an optimal structure that achieves the maximum (or minimum) interested property, *e.g.*, stiffness. Topology structure optimization is crucial to many engineering design and manufacturing problems, including 3D printing, and design of air foils, slab bridges, aerodynamic shapes of race cars, *etc.* The conventional approach is to solve a constraint optimization problem that minimizes a compliance objective subject to a total volume constraint (Sigmund, 1997). However, the numerical computation is usually very costly. We aim to use active learning to learn a model that directly predicts the optimal structure, without the need for running numerical optimization every time.

We considered the stress experiment in (Keshavarzadeh et al., 2018) with an L-shape linear elastic structure. The structure is subjected to a load (*i.e.*, input) on the bottom right half and is discretized in a $[0, 1] \times [0, 1]$ domain. The load is represented by two parameters, the location (in $[0.5, 1]$) and angle (in $[0, \frac{\pi}{2}]$). The goal is to find the structure that achieves the maximum stiffness given the load. Optimizing the structure needs to repeatedly call a numerical solver, where the choice of the mesh determines the fidelity. We used two fidelities to query the training examples. One uses a 50×50 mesh, the other 75×75 . Correspondingly, the output dimensions are 2,500 and 5,625. The costs are measured by the average running time: $\lambda_1 = 1$, $\lambda_2 = 3$. We randomly generated 500 structures for test, where the internal solver uses a 100×100 mesh. Initially, we randomly queried 10 examples at the first fidelity and 2 at the second fidelity. We then ran all the active learning methods to query 100 examples. We conducted the experiments for five times, and report the average nRMSE along with the cost in Fig. 1e. As we can see, DMFAL achieves much better prediction accuracy than the competing approaches (with the same cost). That implies our predicted structures are much closer to the optimal structures. While the performance of the other methods tended to converge early, DMFAL’s

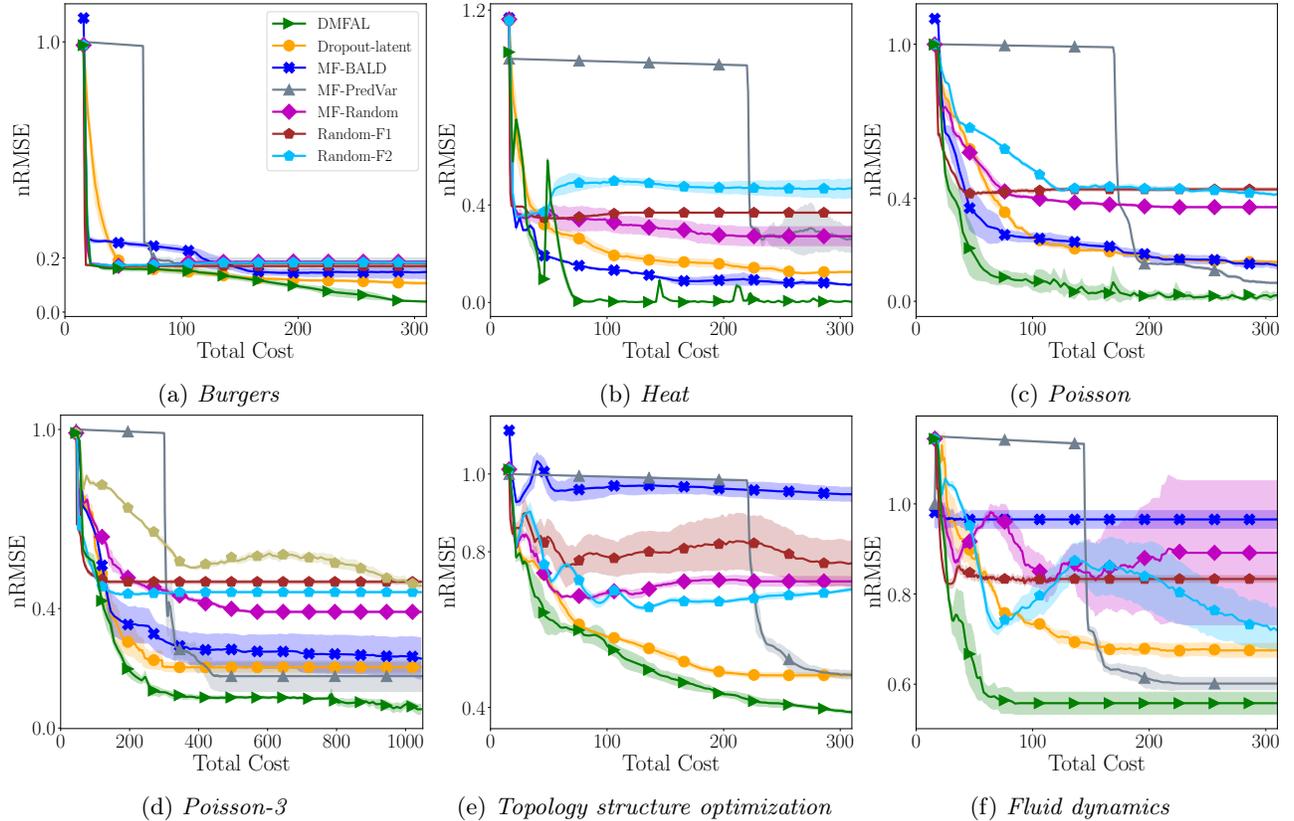


Figure 1: Normalized root-mean-square error (nRMSE) vs. total cost (running time) in active learning. The results are averaged from five runs. The shaded regions indicate the standard deviations.

performance kept improving and the trend did not stop even when all the queries were finished.

For a fine-grained comparison, we visualize eight structures predicted by all the methods, after the active learning is finished. As shown in Fig. 2, DMFAL predicted much more accurate structures, which capture both the global shapes and local details, and the density of the materials is closer to the ground-truth. Although Dropout-latent captures the global shapes as well, its predictions are more blurred and miss many local details, *e.g.*, the second to seventh structure in Fig. 2c. The other methods often provided wrong structures (*e.g.*, the first and last structure in Fig. 2d, f, and g, the second and fourth structure in Fig. 2e) and insufficient density (*e.g.*, the second, third and sixth structure in 2g).

6.3 Predicting Fluid Dynamics

Third, we applied DMFAL in a computational fluid dynamics (CFD) problem. The task is to predict the first component of the velocity field of a flow within a rectangular domain in $[0, 1] \times [0, 1]$. The flow is driven by the boundaries with a prescribed velocity (*i.e.*, input) (Bozeman and Dalton, 1973). Along with time, the local velocities inside the fluid will vary differently,

and eventually result in turbulent flows. Computing these fields along with time requires us to solve the incompressible Navier-Stokes equations (Chorin, 1968), which is known to be challenging to solve because of their complex behaviors under big Reynolds numbers. We considered active learning with two-fidelity queries to predict the first component of the velocity field at evenly spaced 20 time points in $[0, 10]$ (temporal domain). The examples in the first fidelity were generated with a 50×50 mesh in the spatial domain (*i.e.*, $[0, 1] \times [0, 1]$), and the second fidelity 75×75 . The corresponding output dimensions are 50,000 and 112,500. The input is a five dimensional vector that consists of the prescribed boundary velocity and Reynold number. See more details in Sec. A.2 of Appendix. To collect the test dataset, we randomly sampled 256 inputs and computed the solution with a 128×128 mesh. The test outputs are obtained by the cubic-spline interpolation. At the beginning, we randomly queried 10 and 2 training examples in the first and second fidelity, respectively. Then we ran each active learning method with 100 queries. We repeated the experiments for five times. The average nRMSE along with the accumulated cost is shown in Fig. 1f. It can be seen that during the training, DMFAL consistently outperforms all the competing methods by a large margin. On the

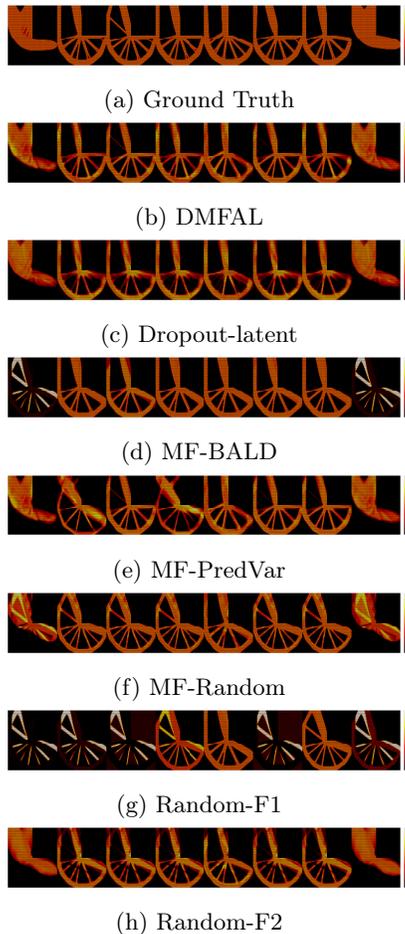


Figure 2: The predicted topology structures for 8 loads. All the active learning approaches started with the same training set (10 fidelity-1 and 2 fidelity-2 examples), and ran with 100 queries.

other hand, to achieve the same level of accuracy, DMFAL spends a much smaller cost. That means, our method requires much less (high-fidelity) simulations to generate the training examples. This is particularly useful for large-scale CFD applications, in which the simulation is known to be very expensive.

To further confirm the gains in computational efficiency, we compared the cost of running the standard numerical methods for the topological optimization and CFD tasks. After active learning, our learned surrogate models give 42x and 466x speed-up in computing (predicting) the high-fidelity solution fields against the numerical methods. Even counting the whole active learning cost, with the growth of test cases, the average cost (time) of DMFAL in computing a solution quickly becomes much smaller than the standard numerical methods. See details in Sec. B of Appendix.

Non-Active Learning. Finally, to confirm the capacity of our deep NN model, we compared with state-of-

the-art high-dimensional non-active learning surrogate models. Note that these models all lack effective active learning strategies. Our model consistently outperforms these methods, often by a large margin (Sec. C of Appendix).

7 CONCLUSION

We have presented DMFAL, a deep multi-fidelity active learning approach for high-dimensional outputs. Our deep neural network based multi-fidelity model is flexibly enough to estimate the strong, complex relationships between the outputs and between the fidelities. We proposed a mutual information based acquisition function that accounts for multi-fidelity queries. To calculate and optimize the acquisition function, we developed an efficient and reliable method that successfully overcomes the computational challenges due to the massive outputs.

Acknowledgments

This work has been supported by MURI AFOSR grant FA9550-20-1-0358 and NSF CAREER Award IIS-2046295.

References

- Ash, J. T., Zhang, C., Krishnamurthy, A., Langford, J., and Agarwal, A. (2019). Deep batch active learning by diverse, uncertain gradient lower bounds. In *International Conference on Learning Representations*.
- Balasubramanian, M. and Schwartz, E. L. (2002). The isomap algorithm and topological stability. *Science*, 295(5552):7–7.
- Balcan, M.-F., Beygelzimer, A., and Langford, J. (2009). Agnostic active learning. *Journal of Computer and System Sciences*, 75(1):78–89.
- Balcan, M.-F., Broder, A., and Zhang, T. (2007). Margin based active learning. In *International Conference on Computational Learning Theory*, pages 35–50. Springer.
- Bickel, P. J. and Doksum, K. A. (2015). *Mathematical statistics: basic ideas and selected topics, volume I*, volume 117. CRC Press.
- Black, F. and Scholes, M. (1973). The pricing of options and corporate liabilities. *Journal of political economy*, 81(3):637–654.
- Bozeman, J. D. and Dalton, C. (1973). Numerical study of viscous flow in a cavity. *Journal of Computational Physics*, 12(3):348–363.
- Burdzy, K., Chen, Z.-Q., Sylvester, J., et al. (2004). The heat equation and reflected brownian motion in

- time-dependent domains. The Annals of Probability, 32(1B):775–804.
- Caretto, L., Gosman, A., Patankar, S., and Spalding, D. (1973). Two calculation procedures for steady, three-dimensional flows with recirculation. In Proceedings of the third international conference on numerical methods in fluid mechanics, pages 60–68. Springer.
- Chapra, S. C., Canale, R. P., et al. (2010). Numerical methods for engineers. Boston: McGraw-Hill Higher Education,.
- Chorin, A. J. (1968). Numerical solution of the navier-stokes equations. Mathematics of computation, 22(104):745–762.
- Chung, T. (2010). Computational fluid dynamics. Cambridge university press.
- Conti, S. and O’Hagan, A. (2010). Bayesian emulation of complex multi-output and dynamic computer models. Journal of statistical planning and inference, 140(3):640–651.
- Cutajar, K., Pullin, M., Damianou, A., Lawrence, N., and González, J. (2019). Deep Gaussian processes for multi-fidelity modeling. arXiv preprint arXiv:1903.07320.
- Damianou, A. and Lawrence, N. (2013). Deep gaussian processes. In Artificial Intelligence and Statistics, pages 207–215.
- Dasgupta, S. (2011). Two faces of active learning. Theoretical computer science, 412(19):1767–1781.
- Ducoffe, M. and Precioso, F. (2018). Adversarial active learning for deep networks: a margin based approach. arXiv preprint arXiv:1802.09841.
- Gal, Y. and Ghahramani, Z. (2016). Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. In international conference on machine learning, pages 1050–1059.
- Gal, Y., Islam, R., and Ghahramani, Z. (2017). Deep Bayesian active learning with image data. In International Conference on Machine Learning, pages 1183–1192.
- Geifman, Y. and El-Yaniv, R. (2017). Deep active learning over the long tail. arXiv preprint arXiv:1711.00941.
- Gissin, D. and Shalev-Shwartz, S. (2019). Discriminative active learning. arXiv preprint arXiv:1907.06347.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial nets. In Advances in neural information processing systems, pages 2672–2680.
- Hamelijnck, O., Damoulas, T., Wang, K., and Girolami, M. (2019). Multi-resolution multi-task gaussian processes. Advances in Neural Information Processing Systems, 32:14025–14035.
- Hanneke, S. et al. (2014). Theory of disagreement-based active learning. Foundations and Trends® in Machine Learning, 7(2-3):131–309.
- Higdon, D., Gattiker, J., Williams, B., and Rightley, M. (2008). Computer model calibration using high-dimensional output. Journal of the American Statistical Association, 103(482):570–583.
- Houlsby, N., Huszár, F., Ghahramani, Z., and Lengyel, M. (2011). Bayesian active learning for classification and preference learning. arXiv preprint arXiv:1112.5745.
- Huang, D., Allen, T. T., Notz, W. I., and Miller, R. A. (2006). Sequential kriging optimization using multiple-fidelity evaluations. Structural and Multidisciplinary Optimization, 32(5):369–382.
- Journel, A. G. and Huijbregts, C. J. (1978). Mining geostatistics, volume 600. Academic press London.
- Kandasamy, K., Dasarathy, G., Oliva, J. B., Schneider, J., and Póczos, B. (2016). Gaussian process bandit optimisation with multi-fidelity evaluations. In Advances in Neural Information Processing Systems, pages 992–1000.
- Kandasamy, K., Dasarathy, G., Schneider, J., and Póczos, B. (2017). Multi-fidelity bayesian optimisation with continuous approximations. In Proceedings of the 34th International Conference on Machine Learning-Volume 70, pages 1799–1808. JMLR. org.
- Kato, T. (2013). Perturbation theory for linear operators, volume 132. Springer Science & Business Media.
- Kennedy, M. C. and O’Hagan, A. (2000). Predicting the output from a complex computer code when fast approximations are available. Biometrika, 87(1):1–13.
- Keshavarzzadeh, V., Kirby, R. M., and Narayan, A. (2018). Parametric topology optimization with multi-resolution finite element models. arXiv preprint arXiv:1808.10367.
- Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980.
- Kingma, D. P. and Welling, M. (2013). Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114.
- Kleinegesse, S. and Gutmann, M. U. (2020). Bayesian experimental design for implicit models by mutual information neural estimation. In International Conference on Machine Learning, pages 5316–5326. PMLR.

- Kutluay, S., Bahadir, A., and Özdeci, A. (1999). Numerical solution of one-dimensional burgers equation: explicit and exact-explicit finite difference methods. Journal of Computational and Applied Mathematics, 103(2):251–261.
- Lam, R., Allaire, D. L., and Willcox, K. E. (2015). Multifidelity optimization using statistical surrogate modeling for non-hierarchical information sources. In 56th AIAA/ASCE/AHS/ASC Structures, Structural Dynamics, and Materials Conference, page 0143.
- Li, S., Xing, W., Kirby, R., and Zhe, S. (2020). Multi-fidelity bayesian optimization via deep neural networks. In Advances in Neural Information Processing Systems.
- McLeod, M., Osborne, M. A., and Roberts, S. J. (2017). Practical bayesian optimization for variable cost objectives. arXiv preprint arXiv:1703.04335.
- Mockus, J. (2012). Bayesian approach to global optimization: theory and applications, volume 37. Springer Science & Business Media.
- Nagel, K. (1996). Particle hopping models and traffic flow theory. Physical review E, 53(5):4655.
- Oehlert, G. W. (1992). A note on the delta method. The American Statistician, 46(1):27–29.
- Olsen-Kettle, L. (2011). Numerical solution of partial differential equations. Lecture notes at University of Queensland, Australia.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al. (2019). Pytorch: An imperative style, high-performance deep learning library. In Advances in neural information processing systems, pages 8026–8037.
- Perdikaris, P., Raissi, M., Damianou, A., Lawrence, N., and Karniadakis, G. E. (2017). Nonlinear information fusion algorithms for data-efficient multi-fidelity modelling. Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences, 473(2198):20160751.
- Persides, S. (1973). The laplace and poisson equations in schwarzschild’s space-time. Journal of Mathematical Analysis and Applications, 43(3):571–578.
- Picheny, V., Ginsbourger, D., Richet, Y., and Caplin, G. (2013). Quantile-based optimization of noisy computer experiments with tunable precision. Technometrics, 55(1):2–13.
- Poloczek, M., Wang, J., and Frazier, P. (2017). Multi-information source optimization. In Advances in Neural Information Processing Systems, pages 4288–4298.
- Raissi, M., Perdikaris, P., and Karniadakis, G. E. (2017). Physics informed deep learning (part i): Data-driven solutions of nonlinear partial differential equations. arXiv preprint arXiv:1711.10561.
- Raissi, M., Perdikaris, P., and Karniadakis, G. E. (2019). Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. Journal of Computational Physics, 378:686–707.
- Rasmussen, C. E. and Ghahramani, Z. (2002). Infinite mixtures of gaussian process experts. Advances in neural information processing systems, 2:881–888.
- Rozvany, G. I. (2009). A critical review of established methods of structural topology optimization. Structural and multidisciplinary optimization, 37(3):217–237.
- Sacks, J., Welch, W. J., Mitchell, T. J., and Wynn, H. P. (1989). Design and analysis of computer experiments. Statistical science, 4(4):409–423.
- Schölkopf, B., Smola, A., and Müller, K.-R. (1998). Nonlinear component analysis as a kernel eigenvalue problem. Neural computation, 10(5):1299–1319.
- Sener, O. and Savarese, S. (2018). Active learning for convolutional neural networks: A core-set approach. In International Conference on Learning Representations.
- Settles, B. (2009). Active learning literature survey. Technical report, University of Wisconsin-Madison Department of Computer Sciences.
- Settles, B., Craven, M., and Friedland, L. (2008). Active learning with real annotation costs. In Proceedings of the NIPS workshop on cost-sensitive learning, volume 1. Vancouver, CA:.
- Shah, A., Xing, W., and Triantafyllidis, V. (2017). Reduced-order modelling of parameter-dependent, linear and nonlinear dynamic partial differential equation models. Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences, 473(2200):20160809.
- Sigmund, O. (1997). On the design of compliant mechanisms using topology optimization. Journal of Structural Mechanics, 25(4):493–524.
- Snoek, J., Larochelle, H., and Adams, R. P. (2012). Practical bayesian optimization of machine learning algorithms. In Advances in neural information processing systems, pages 2951–2959.
- Snoek, J., Rippel, O., Swersky, K., Kiros, R., Satish, N., Sundaram, N., Patwary, M., Prabhat, M., and Adams, R. (2015). Scalable bayesian optimization using deep neural networks. In International conference on machine learning, pages 2171–2180.

- Song, J., Chen, Y., and Yue, Y. (2019). A general framework for multi-fidelity bayesian optimization with gaussian processes. In The 22nd International Conference on Artificial Intelligence and Statistics, pages 3158–3167.
- Spitzer, F. (1964). Electrostatic capacity, heat flow, and brownian motion. Probability theory and related fields, 3(2):110–121.
- Sugimoto, N. (1991). Burgers equation with a fractional derivative; hereditary effects on nonlinear acoustic waves. Journal of fluid mechanics, 225:631–653.
- Takeno, S., Fukuoka, H., Tsukada, Y., Koyama, T., Shiga, M., Takeuchi, I., and Karasuyama, M. (2019). Multi-fidelity bayesian optimization with max-value entropy search. arXiv preprint arXiv:1901.08275.
- Versteeg, H. K. and Malalasekera, W. (2007). An introduction to computational fluid dynamics: the finite volume method. Pearson education.
- Wainwright, M. J., Jordan, M. I., et al. (2008). Graphical models, exponential families, and variational inference. Foundations and Trends® in Machine Learning, 1(1-2):1–305.
- Wang, Z. and Jegelka, S. (2017). Max-value entropy search for efficient bayesian optimization. In Proceedings of the 34th International Conference on Machine Learning-Volume 70, pages 3627–3635. JMLR. org.
- Wang, Z., Xing, W., Kirby, R., and Zhe, S. (2021). Multi-fidelity high-order gaussian processes for physical simulation. In International Conference on Artificial Intelligence and Statistics, pages 847–855. PMLR.
- Wilson, A. G., Knowles, D. A., and Ghahramani, Z. (2012). Gaussian process regression networks. In Proceedings of the 29th International Conference on Machine Learning, pages 1139–1146.
- Wu, J. and Frazier, P. I. (2017). Continuous-fidelity bayesian optimization with knowledge gradient. In NIPS Workshop on Bayesian Optimization.
- Xing, W., Shah, A. A., and Nair, P. B. (2015). Reduced dimensional gaussian process emulators of parametrized partial differential equations based on isomap. In Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences, volume 471, page 20140697. The Royal Society.
- Xing, W., Triantafyllidis, V., Shah, A., Nair, P., and Zabaras, N. (2016). Manifold learning for the emulation of spatial fields from computational models. Journal of Computational Physics, 326:666–690.
- Zhe, S., Xing, W., and Kirby, R. M. (2019). Scalable high-order gaussian process regression. In The 22nd International Conference on Artificial Intelligence and Statistics, pages 2611–2620.
- Zienkiewicz, O. C., Taylor, R. L., Zienkiewicz, O. C., and Taylor, R. L. (1977). The finite element method, volume 36. McGraw-hill London.

Supplementary Material: Deep Multi-Fidelity Active Learning of High-Dimensional Outputs

A EXPERIMENTAL DETAILS

A.1 Solving Partial Differential Equations

Burgers’ equation is a canonical nonlinear hyperbolic PDE, widely used to model various physical phenomena, such as nonlinear acoustics (Sugimoto, 1991), fluid dynamics (Chung, 2010), and traffic flows (Nagel, 1996). Due to its capability of developing discontinuities (*i.e.*, shock waves), Burger’s equation is used as a benchmark test example for many numerical solvers and surrogate models (Kutluay et al., 1999; Shah et al., 2017; Raissi et al., 2017). The viscous version of Burger’s equation is given by

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} = v \frac{\partial^2 u}{\partial x^2},$$

where u is the volume, x is a spatial location, t is the time, and v is the viscosity. We set $x \in [0, 1]$, $t \in [0, 3]$, and $u(x, 0) = \sin(x\pi/2)$ with a homogeneous Dirichlet boundary condition. The input parameter is the viscosity $v \in [0.001, 0.1]$. Given the input, we aim to predict the solution field (*i.e.*, the values of u) in the spatial-temporal domain $[0, 1] \times [0, 3]$. To obtain the training and test datasets, we solve the equation using the finite element (Zienkiewicz et al., 1977) with hat functions in space and backward Euler in time domains on a regular mesh.

Poisson’s equation is an elliptic PDE and commonly used to model potential fields, *e.g.*, electrostatic and gravitational fields, in physics and mechanical engineering (Chapra et al., 2010). The equation used in our experiment is given by

$$\Delta u = \beta \delta(\mathbf{x} - \mathbf{c}),$$

where Δ is the Laplace operator (Persides, 1973), \mathbf{u} is the volume, $\delta(\cdot)$ is the Dirac-delta function, and \mathbf{c} is the center of the domain. We used a 2D spatial domain, $\mathbf{x} \in [0, 1] \times [0, 1]$, and Dirichlet boundary conditions. We used the constant values of the four boundaries and β as the input parameters, each of which ranges from 0.1 to 0.9. We solved the equation using the finite difference method with the first order center differencing scheme and regular rectangle meshes.

Heat equation is a fundamental PDE that models heat conduction over time. It is also widely used in many other areas, such as probability theory (Spitzer, 1964; Burdzy et al., 2004) and financial mathematics (Black and Scholes, 1973). The equation is defined as

$$\frac{\partial u}{\partial t} + \alpha \Delta u = 0,$$

where u is the heat, α the thermal conductivity, and Δ the Laplace operator. In our experiment, we used a 2D spatial-temporal domain $x \in [0, 1]$, $t \in [0, 5]$ with the Neumann boundary condition at $x = 0$ and $x = 1$, and $u(x, 0) = H(x - 0.25) - H(x - 0.75)$, where $H(\cdot)$ is the Heaviside step function. We considered three input parameters — the flux rate $\in [0, 1]$ of the left boundary at $x = 0$, the flux rate $\in [-1, 0]$ of the right boundary at $x = 1$, and $\alpha \in [0.01, 0.1]$. To generate the training and test examples, we solve the equation with the finite difference in the space domain and backward Euler in the time domain.

A.2 Predicting Fluid Dynamics

We also examined DMFAL in predicting the velocity field of a flow within a rectangular domain with a prescribed velocity along the boundaries (Bozeman and Dalton, 1973). This is a classical computational fluid dynamics (CFD) problem. The simulation of the flow involves solving the incompressible Navier-Stokes (NS) equation (Chorin, 1968),

$$\rho(\mathbf{u} \cdot \nabla) \mathbf{u} = -\nabla p + \mu \nabla^2 \mathbf{u},$$

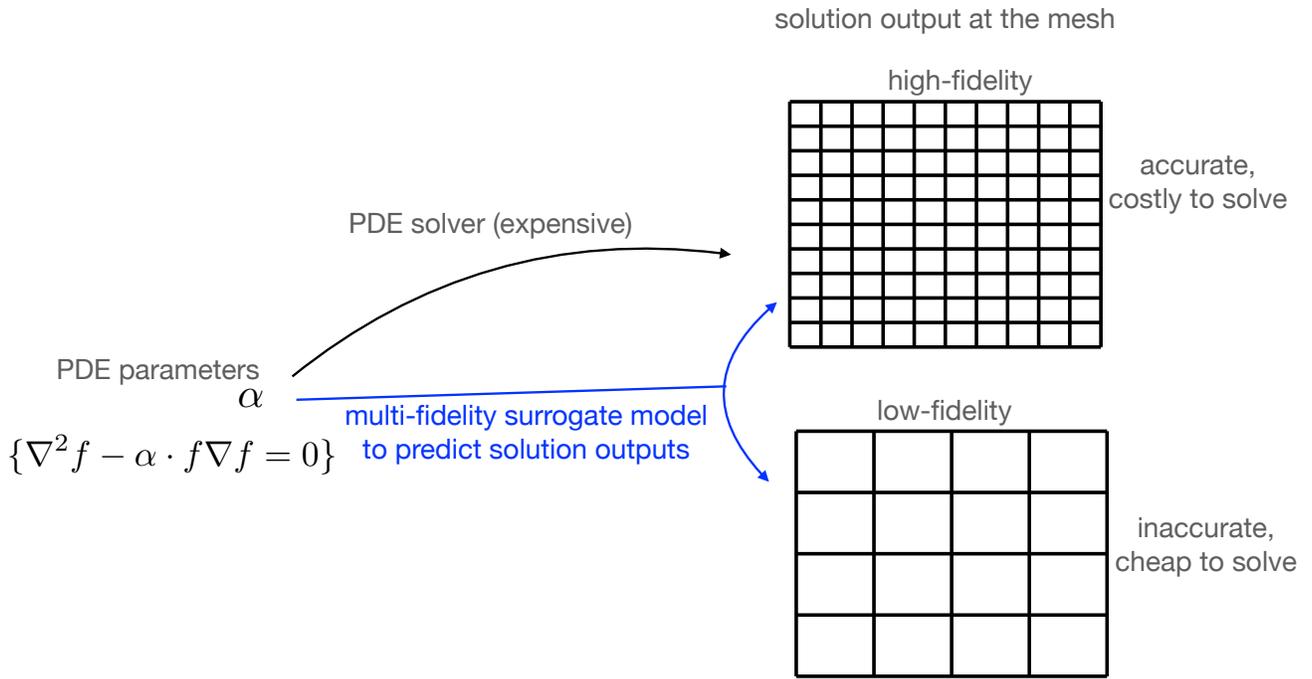


Figure 3: Illustration of the motivation and goal with physical simulation as an example. Solving every PDE from scratch is expensive. Hence, we learn a multi-fidelity surrogate model that can predict high-fidelity solution outputs outright given the PDE parameters. We develop active learning methods to further reduce the cost of running numerical solvers to generate/collect training examples.

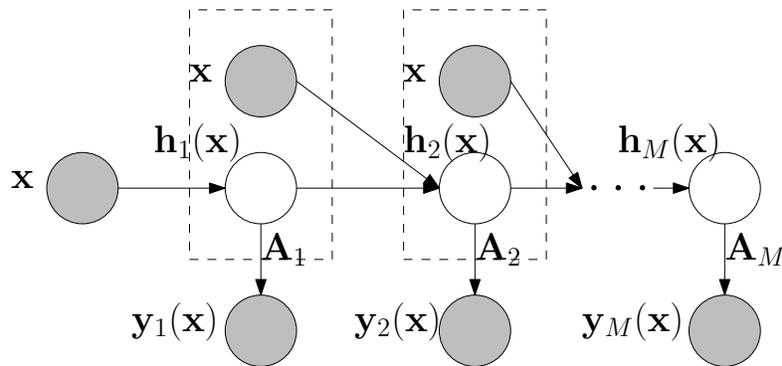


Figure 4: Graphical representation of the deep multi-fidelity model. The low dimensional latent output in each fidelity $\mathbf{h}_m(\mathbf{x})$ ($1 \leq m \leq M$) is generated by a (deep) neural network.

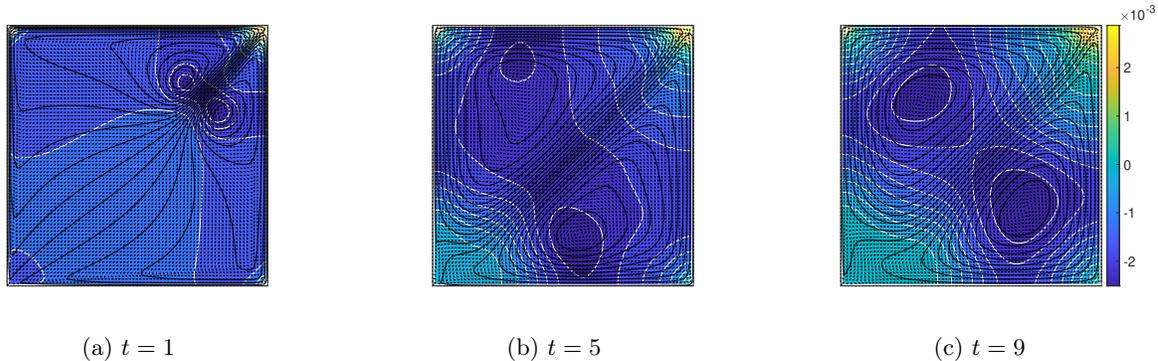


Figure 5: Examples of the first component of the velocity field (with contour lines) at three time points.

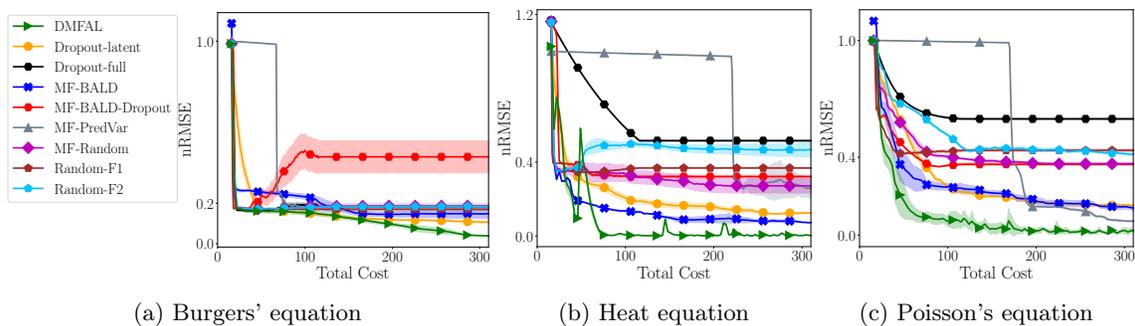


Figure 6: Normalized root-mean-square error (nRMSE) for active learning of PDE solution fields with two-fidelity queries. The normalizer is the mean of the test outputs. The results are averaged from five runs. The shaded regions indicate the standard deviations.

where ρ is the density, p is the pressure, \mathbf{u} is the velocity, and μ is the dynamic viscosity. The equation is well known to be challenging to solve due to their complicated behaviours under large Reynolds numbers. We set the rectangular domain to $[0, 1] \times [0, 1]$, and time $t \in [0, 10]$. The input includes the tangential velocities of the four boundaries and the Reynolds number $\in [100, 5000]$. The output are the first component of the velocity field at 20 equally spaced time steps in $[0, 10]$ (see Fig. 5). To generate the training and test examples, we used the SIMPLE algorithm (Caretto et al., 1973) with a stagger grid (Versteeg and Malalasekera, 2007), the up-wind scheme (Versteeg and Malalasekera, 2007) for the spatial difference, and the implicit time scheme with fixed time steps to solve the NS equation.

A.3 Using Full Dropout

In the experiment of Section 6.1 of the main paper, we have also applied MC dropout to directly generate the posterior samples of the final output in each fidelity $\{\mathbf{y}_m\}_{m=1}^M$. We generated 100 samples. We used these high-dimensional samples to calculate the empirical means and covariance matrices, based on which we estimated a multi-variate Gaussian posterior for each \mathbf{y}_m and $\hat{\mathbf{y}}_m = [\mathbf{y}_m, \mathbf{y}_M]$ via moment matching. These posterior distributions are then used to calculate and optimize our acquisition function and multi-fidelity BALD in the active learning. We denote these methods by Dropout-full and MF-BALD-Dropout. We report their nRMSE *vs.* cost in Fig. 6 and 7, along with all the other methods. As we can see, in all the cases, Dropout-full is far worse than Dropout-latent and even inferior to random query strategies, except that in Fig. 6a, they are quite close. The prediction accuracy of MF-BALD is always much better than MF-BALD-Dropout. Those results demonstrate that the posterior distributions of the high-dimensional outputs, if fully estimated by a small number of dropout samples, are far less accurate and reliable than our method. Also, the computation is much more costly — we need to directly compute an empirical covariance matrix based on these high-dimensional samples, and calculate the log determinant.

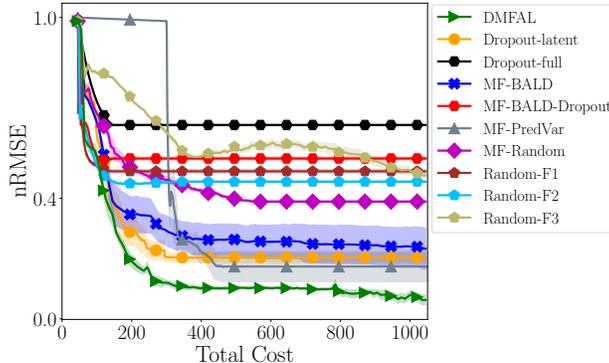


Figure 7: Normalized root-mean-square error (nRMSE) for active learning of the solution field of Poisson’s equation, with three-fidelity queries. The results are averaged from five runs.

B GAINS IN COMPUTATIONAL EFFICIENCY

To confirm the gains in computational efficiency, we compared the cost of running standard numerical methods for topology structure optimization and fluid dynamics (Sec 6.2 and 6.3 in the main paper). First, after the active learning, we examined the average time of our learned model in computing (or predicting) a high-fidelity solution field, and contrast to that of the standard numerical approaches. As shown in Fig. 8a, our model is much faster, giving 42x and 466x speed-up for the two tasks. Next, we consider adding the cost (running time) of the active learning, and calculate the average cost of predicting a high-fidelity solution field,

$$\text{Cost}_{\text{avg}} = \frac{\text{Active Learning Cost} + \text{Prediction Cost}}{\text{Number of Acquired Solutions}}.$$

Note that for the numerical methods, there is no learning procedures and the active learning cost is zero. We examined how the average cost varies along with more acquired solutions (*i.e.*, test cases). As we can see in Fig. 8 b and c, when only a few solutions are needed, directly running numerical methods is more economic because we do not need to conduct an extra learning procedure. However, when we need to compute more and more solutions (test cases), the average cost of using the surrogate model quickly becomes much smaller, and keeps decreasing. That is because after training, the surrogate model can predict solutions way more efficiently than running numerical methods from scratch, and the training cost will be diluted by the number of solutions acquired, finally becoming negligible. All these results have demonstrated the gains of computational efficiency of our method.

C NON-ACTIVE LEARNING

To confirm the surrogate performance of our deep NN model, we also compared with several state-of-the-art high-dimensional non-active learning surrogate models: (1) PCA-GP (Higdon et al., 2008), (2) IsoMap-GP (Xing et al., 2015), (3) KPCA-GP (Xing et al., 2016), and (4) HOGP (Zhe et al., 2019). PCA-GP, IsoMap-GP and KPCA-GP are based on the classical linear model of coregionalization (LMC) (Journel and Huijbregts, 1978), and obtain the bases or low-rank structures from Principal Component Analysis (PCA), IsoMap (Balasubramanian and Schwartz, 2002) and Kernel PCA (Schölkopf et al., 1998), respectively. HOGP tensorizes the outputs, generalizes a multilinear Bayesian model with the kernel trick, and is flexible enough to capture nonlinear output correlations and can efficiently deal with very high-dimensional outputs. These methods are all single-fidelity models. In addition, we also compare with (6) MFHoGP (Wang et al., 2021), a state-of-the-art multi-fidelity high-dimensional surrogate model based on matrix GPs. Note that MFHoGP does not support varying output dimensions across the fidelities.

We tested all these methods in the three applications: *Poisson-3*, *Topology structure optimization* and *Fluid dynamics* (see Sec. 6.1, 6.2 and 6.3 of the main paper). For *Poisson-3*, we randomly generated 43, 43, 42 examples for the first, second and third fidelity, respectively. The inputs were uniformly sampled from the input space. For *Topology structure optimization* and *Fluid dynamics*, we randomly generated 64 examples for each fidelity (two fidelities in total). For each application, we randomly generated 512 examples for testing. The settings of those

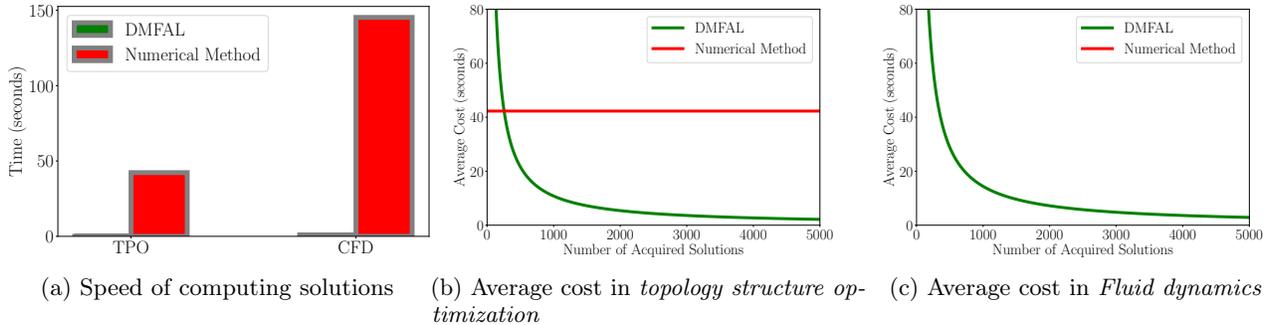


Figure 8: Speed of computing solutions after active learning (a) and total average solving (or predicting) cost with active learning included.

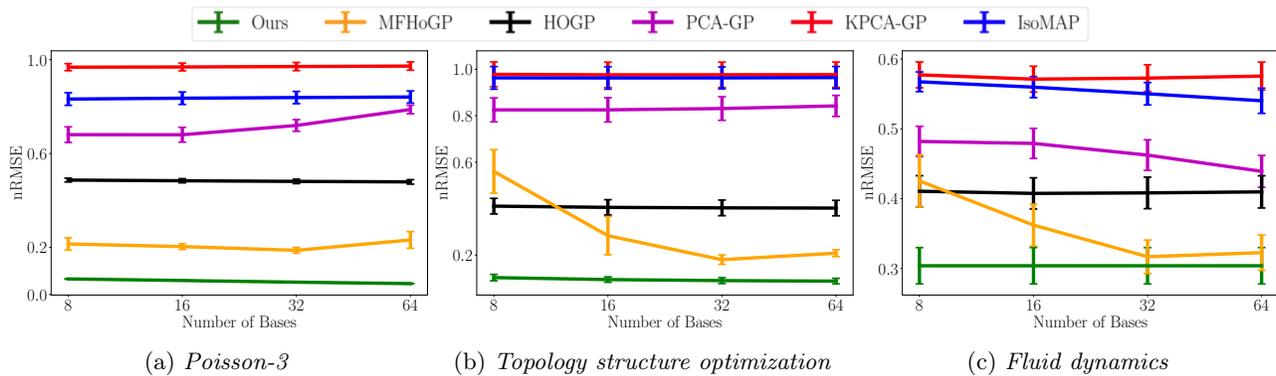


Figure 9: Predictive performance of non-active learning. The results were averaged over five runs.

fidelities are the same as in the main paper. Since MFHoGP does not support varying output dimensions in different fidelities, we interpolated the training outputs at low fidelities to ensure all the fidelities have the same output dimension. We varied the number of bases (required by all the competing methods) from $\{8, 16, 32, 64\}$, which corresponds to the dimension of the latent output in our model (see Sec. 3 of the main paper). We ran the experiments for five times, and report the average nRMSE and its standard deviation in Fig. 9. As we can see, our model consistently outperforms all the methods, often by a large margin. MFHoGP is the second best, but in most cases, its prediction accuracy is still significantly worse than our model ($p < 0.05$). The remaining single-fidelity models are almost always worse than the multi-fidelity models — that is reasonable, because the former cannot differentiate examples of different fidelities, blindly mix them in the training, and hence cannot leverage their relationships to improve the prediction. Together these results have shown the advantage of our deep surrogate model, even in the non-active learning setting.

D ABLATION STUDY ABOUT MULTI-VARIATE DELTA METHOD

We examined the rationality of the multi-variate delta method used in our work (see (6) in Section 4.2 of the main paper). To this end, we used a Bayesian neural networks (BNN) to learn two nontrivial benchmark functions, *Branin* and *Levy*, which are defined as follows. For *Branin*, the input $\mathbf{x} \in [-5, 10] \times [0, 15]$, and for *Levy*, $\mathbf{x} \in [-10, 10]^2$. We used three hidden layers, with 40 neurons per layer, and \tanh as the activation function (the same as in our experiment).

$$\begin{aligned}
 y_{\text{Branin}}(\mathbf{x}) &= - \left(\frac{-1.275x_1^2}{\pi^2} + \frac{5x_1}{\pi} + x_2 - 6 \right)^2 - \left(10 - \frac{5}{4\pi} \right) \cos(x_1) - 10, \\
 y_{\text{Levy}}(\mathbf{x}) &= - \sin^2(3\pi x_1) - (x_1 - 1)^2 [1 + \sin^2(3\pi x_2)] - (x_2 - 1)^2 [1 + \sin^2(2\pi x_2)].
 \end{aligned} \tag{13}$$

We generated N training examples by randomly sampling from the input domain, and then used the structural variational inference to estimate the BNN. In our work, we used the first-order Taylor expansion of the NN output

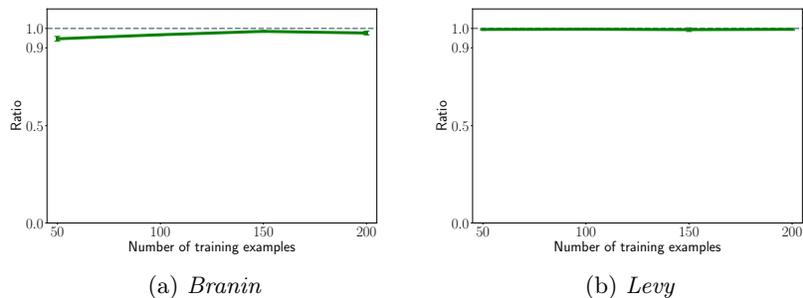


Figure 10: The ratio between first-order and second-order approximations.

to approximate the posterior mean and (co-)variance (see (6-8) in the main paper). To confirm the rationality, we examined the ratio between the first-order Taylor expansion and the second-order Taylor expansion. Specifically, given a new input \mathbf{x} , we first sample an instance of the NN weights \mathcal{W} from the learned posterior $p(\mathcal{W}|\mathcal{D})$. Then, we expand the NN output $f_{\mathcal{W}}(\mathbf{x})$ at the posterior mean of the weights $\mathbb{E}[\mathcal{W}]$, and calculate

$$\text{Ratio} = \frac{|\text{First-Order Taylor Expansion of } f_{\mathcal{W}}(\mathbf{x})|}{|\text{Second-Order Taylor Expansion of } f_{\mathcal{W}}(\mathbf{x})|}.$$

We varied N from $\{50, 100, 150, 200\}$. To obtain a reliable estimate of the ratio, for each N , we randomly sampled 100 inputs. Given each input, we randomly drew 10 samples of \mathcal{W} to calculate the ratio. We report the average ratio, and its standard deviation in Fig. 10. As we can see, in both cases, the ratio is constantly close to one, and the standard deviation is close to zero. That means, the first-order approximation dominates, and the second-order terms can be ignored. This is consistent with our analysis in the main paper.