

# Fast Regression of the Tritium Breeding Ratio in Fusion Reactors

P Mánek<sup>1,2</sup>, G Van Goffrier<sup>1</sup>, V Gopakumar<sup>3</sup>, N Nikolaou<sup>1</sup>,  
J Shimwell<sup>3</sup> and I Waldmann<sup>1</sup>

<sup>1</sup> Department of Physics and Astronomy, University College London, Gower Street, London WC1E 6BT, UK

<sup>2</sup> Institute of Experimental and Applied Physics, Czech Technical University, Husova 240/5, Prague 110 00, Czech Republic

<sup>3</sup> UK Atomic Energy Authority, Culham Science Centre, OX14 3DB Abingdon, UK

E-mail: petr.manek.19@ucl.ac.uk, graham.vangoffrier.19@ucl.ac.uk

**Abstract.** The tritium breeding ratio (TBR) is an essential quantity for the design of modern and next-generation D-T fueled nuclear fusion reactors. Representing the ratio between tritium fuel generated in breeding blankets and fuel consumed during reactor runtime, the TBR depends on reactor geometry and material properties in a complex manner. In this work, we explored the training of surrogate models to produce a cheap but high-quality approximation for a Monte Carlo TBR model in use at the UK Atomic Energy Authority. We investigated possibilities for dimensional reduction of its feature space, reviewed 9 families of surrogate models for potential applicability, and performed hyperparameter optimisation. Here we present the performance and scaling properties of these models, the fastest of which, an artificial neural network, demonstrated  $R^2 = 0.985$  and a mean prediction time of  $0.898 \mu\text{s}$ , representing a relative speedup of  $8 \cdot 10^6$  with respect to the expensive MC model. We further present a novel adaptive sampling algorithm, Quality-Adaptive Surrogate Sampling, capable of interfacing with any of the individually studied surrogates. Our preliminary testing on a toy TBR theory has demonstrated the efficacy of this algorithm for accelerating the surrogate modelling process.

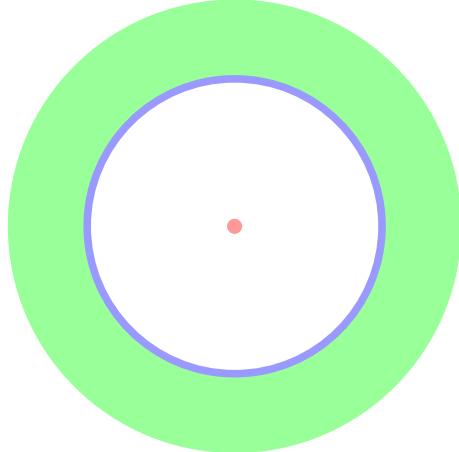
**Index terms**— Machine Learning, Surrogate Model, Tritium Breeding, Regression, Fast Approximation, Adaptive Sampling Submitted to: *Nucl. Fusion*

## 1. Introduction

Surrogate models were developed to resolve computational limitations in the analysis of massive datasets by replacing a resource-expensive procedure with a much cheaper approximation [?]. They are especially useful in applications where numerous evaluations of an expensive procedure are required over the same or similar domains, e.g. in the parameter optimisation of a theoretical model. The term “metamodel” proves especially meaningful in this case, when the surrogate model approximates a computational process which is itself a model for a (perhaps unknown) physical process [?]. There exists a spectrum between “physical” surrogates which are constructed with some contextual knowledge in hand, and “empirical” surrogates which are derived purely from samples of the underlying expensive model.

In this work, we develop a family of empirical surrogate models for the tritium breeding ratio (TBR) in an inertial confinement fusion (ICF) reactor. The expensive model that our surrogate model approximates is a Monte Carlo (MC) neutronics simulation, Paramak [?], which returns a prediction of the TBR for a given configuration of a spherical ICF reactor. Although more expensive 3D parametric models exist, we chose the Paramak simulation for its preferable speed in dataset generation in order to most fully demonstrate our methods. We quantify the success of several of our best-performing surrogate models by studying their accuracy and prediction time. We further propose an adaptive sampling algorithm (QASS) suitable for reducing the quantity of expensive samples needed to train our surrogate models.

Paramak facilitates simulation via an OpenMC neutronics workflow that is en-



**Figure 1.** Diagram of the simple sphere geometry (not to scale) where the blanket is **green**, the first wall is **blue** and the neutron point source is **red**. Blanket and first wall thickness, as well as their material and structural properties, are adjustable parameters of the simulation that are later optimized (see Table 1 for complete parameter listing).

closed in a portable Docker container, which conveniently exposes an HTTP API using the Python 3 `flask` package. Within this setup, we employ a Muir energy distribution [?]† around 14.06 MeV to approximate a Deuterium-Tritium (D-T) plasma neutron source. Illustrated in Figure 1, the simulated reactor geometry was left adjustable in order to study its influence on the TBR. Nuclear data for simulation was extracted from the following sources, in order of preference: FENDL 3.1d [?]; JEFF 3.3 [?]; and ENDF/B-VII.1 [?]. To maintain model-agnostic approach, variance reduction (VR) techniques were not used to accelerate the MC neutronics simulation [?]. It should be noted that depending on application, VR may constitute a viable alternative to the presented work.

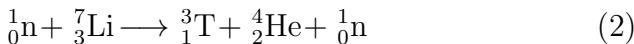
For the remainder of Section 1, we will define the TBR and further motivate our research. In Section 2 we will present our methodologies for the comparison testing of a wide variety of surrogate modelling techniques,

as well as defining the adaptive sampling procedure QASS. After delivering the results of these approaches in Section 3, we will give our final conclusions and recommendations in Section 4.

### 1.1. Problem Description

Nuclear fusion technology relies on the production and containment of an extremely hot and dense plasma containing enriched Hydrogen isotopes. The current frontier generation of fusion reactors, such as the Joint European Torus (JET) and the under-construction **ITER**, make use of both tritium and deuterium fuel. While at least one deuterium atom occurs for every 5000 molecules of naturally-sourced water, and may be easily distilled, tritium is extremely rare in nature. It may be produced indirectly through irradiation of heavy water ( $D_2O$ ) during nuclear fission, but only at very low rates which could never sustain industrial-scale fusion power.

**Modern** D-T reactors rely on tritium breeding blankets, specialized layers of material which partially line the reactor and produce tritium upon neutron bombardment, e.g. by



where T represents tritium and  $^7Li$ ,  $^6Li$  are the more and less common isotopes of lithium, respectively. The TBR is defined as the ratio of tritium produced per source neutron , whose description in Paramak is facilitated by 2 classes of parameters (exhaustively listed in Table 1). While the geometry of a given reactor is described by continuous parameters, material selections are specified by discrete

† A bug in the Muir distribution involving erroneous normal sampling was recently uncovered (see <https://github.com/openmc-dev/openmc/pull/1670>), but is disregarded in the present approach-focused work.

**Table 1.** Input parameters supplied to Paramak and surrogates in alphabetical order. Groups of fractions marked<sup>†‡</sup> are independently required to sum to 1.

	Parameter name (abbreviation)	Domain
Blanket	Breeder fraction <sup>†</sup>	[0,1]
	Breeder $^6Li$ enrichment fraction	[0,1]
	Breeder material (BBM)	{ $Li_2TiO_3, Li_4SiO_4$ }
	Breeder packing fraction	[0,1]
	Coolant fraction <sup>†</sup>	[0,1]
	Coolant material (BCM)	{ $D_2O, H_2O, He$ }
	Multiplier fraction <sup>†</sup>	[0,1]
	Multiplier material (BMM)	{Be, $Be_{12}Ti$ }
	Multiplier packing fraction	[0,1]
	Structural fraction <sup>†</sup> (BSM)	[0,1]
First wall	Structural material	{SiC, eurofer}
	Thickness	[0,500] cm
	Armour fraction <sup>‡</sup>	[0,1]
	Coolant fraction <sup>‡</sup>	[0,1]
	Coolant material (FCM)	{ $D_2O, H_2O, He$ }
Structural	Structural fraction <sup>‡</sup>	[0,1]
	Structural material (FSM)	{SiC, eurofer}
	Thickness	[0,20] cm

categorical parameters. For all parameters, we have attempted to cover the full theoretical range of values even where those values are practically infeasible with current technology (e.g. packing fractions close to 1). Simulating broadly around typical values of parameters is necessary to construct the most accurate model of the typical range, and also aids in demonstrating the robustness of constructed models.

In our work, we set out to produce a fast TBR surrogate model, which takes the same input parameters as the MC model used in Paramak and approximates its output with the greatest achievable regression performance, while also minimising the required quantity of expensive-model samples needed for training.

**Table 2.** Considered surrogate model families, their selected abbreviations and implementations.  $\mathcal{H}$  denotes the set of hyperparameters, family-dependent priors that control the learning process, and are tuned separately. Families with fewer hyperparameter represents a smaller surrogate domain to explore.

Surrogate family	Abbr.	Impl.	$ \mathcal{H} $
Support vector machines [?]	SVM	SciKit [?]	3
Gradient boosted trees [?, ?, ?]	GBT	SciKit	11
Extremely randomized trees [?]	ERT	SciKit	7
AdaBoosted decision trees <sup>a</sup> [?]	ABT	SciKit	3
Gaussian process regression [?]	GPR	SciKit	2
$k$ nearest neighbours	KNN	SciKit	3
Artificial neural networks	ANN	Keras [?]	2
Inverse distance weighing [?]	IDW	SMT [?]	1
Radial basis functions	RBF	SMT	3

<sup>a</sup>Note that ABTs can be viewed as a subclass of GBTs.

## 2. Methodology

Labeling the expensive Paramak model  $f(x)$ , a surrogate is a mapping  $\hat{f}(x)$  such that  $f(x)$  and  $\hat{f}(x)$  minimize a selected dissimilarity metric. In order to be considered *viable*,  $\hat{f}(x)$  is required to achieve expected evaluation time lower than that of  $f(x)$ . In this work, we consider 2 methods of producing viable surrogates: (1) a conventional decoupled approach, which evaluates  $f(x)$  on a set of randomly sampled points and trains surrogates in a supervised scheme, and (2) an adaptive approach, which attempts to compensate for localized regression performance insufficiency's by interleaving multiple epochs of sampling and training.

For both methods, we selected state-of-the-art regression algorithms to perform surrogate training on sampled point sets. Listed in Table 2, these implementations define 9 surrogate families that are later reviewed in Section 3. We note that each presented algorithm defines hyperparameters that may influence its performance. Their problem-specific optimal values are searched

within the scope of this work, in particular in Experiments 1 & 2 that are outlined in Section 2.1.

To compare the quality of the produced surrogates, we define a number of metrics listed in Table 3. For regression performance analysis, we include a selection of absolute metrics (**MAE**, **S**) to assess the models' approximation capability and set practical bounds on the expected uncertainty of their predictions. In addition, we also track relative measures ( **$R^2$** ,  **$R_{adj}^2$** ) that are better-suited for comparison between this work and others as they are invariant with respect to the selected domain and image space. For analysis of computational complexity, surrogates are assessed in terms of wall time (captured by the Python 3 `time` package). This is motivated by common practical use cases of our work, where models are trained and used as drop-in replacements for Paramak. All times reported (training, test, evaluation) are normalized by the corresponding dataset size, i.e. correspond to "time to process a single datapoint."

Even though some surrogates support acceleration by means of **parallelization**, we used non-parallelized implementations. The only exception to this is the ANN family, which requires a considerable amount of processing power for training on conventional CPU architectures. Lastly, to prevent undesirable bias by training set selection, all reported metrics are obtained via 5-fold cross-validation. In this setting, a sample set is uniformly divided into 5 disjoint folds, each of which is used as a test set for models trained on the remaining 4. Having repeated the same experiment for each division, the overall value of individual metrics is reported in terms of their mean and standard deviation over all folds.

**Table 3.** Metrics recorded in experiments. In formulations, we work with a training set of size  $N_0$  and a test set of size  $N$ , values  $y^{(i)} = f(x^{(i)})$  and  $\hat{y}^{(i)} = \hat{f}(x^{(i)})$  denote images of the  $i$ th testing sample in Paramak and the surrogate respectively. The mean  $\bar{y} = N^{-1} \sum_{i=1}^N y^{(i)}$  and  $P$  is the number of input features.

Regression performance metrics	Notation	Mathematical formulation
Mean absolute error	MAE	$N^{-1} \sum_{i=1}^N  y^{(i)} - \hat{y}^{(i)} $
Standard deviation of error	$S$	$\text{StdDev}_{i=1}^N \{ y^{(i)} - \hat{y}^{(i)} \}$
Coefficient of determination	$R^2$	$1 - \sum_{i=1}^N (y^{(i)} - \hat{y}^{(i)})^2 \left[ \sum_{i=1}^N (y^{(i)} - \bar{y})^2 \right]^{-1}$
Adjusted $R^2$	$R_{\text{adj.}}^2$	$1 - (1 - R^2)(N - 1)(N - P - 1)^{-1}$
Computational complexity metrics		
Mean training time	$\bar{t}_{\text{trn.}}$	(wall training time of $\hat{f}(x)$ ) $N_0^{-1}$
Mean prediction time	$\bar{t}_{\text{pred.}}$	(wall prediction time of $\hat{f}(x)$ ) $N^{-1}$
Relative speedup	$\omega$	(wall evaluation <sup>b</sup> time of $f(x)$ ) $(N \bar{t}_{\text{pred.}})^{-1}$

<sup>b</sup>This corresponds to evaluation of Paramak on all points of the test set. In surrogates, the equivalent time period is referred to as “prediction time.”

## 2.1. Decoupled Approach

Experiments related to the decoupled approach are organized in 4 parts, further described in this section. In summary, we aim to optimize the hyperparameters of each surrogate family separately and later use the best found results to compare surrogate families among themselves.

The objective of Experiment 1 is to simplify the regression task for surrogates prone to sub-optimal performance in discontinuous spaces. To this end, training points are filtered to a single selected discrete feature assignment, and surrogates are trained only on the remaining continuous features. This is repeated several times to explore variances in behaviour, particularly in 4 distinct assignments that are obtained by setting blanket and first wall coolant materials to one of: {H<sub>2</sub>O, He}. Experiment 2 conventionally measures surrogate performance on the full feature space without any parameter restrictions. In both experiments, hyperparameter tuning is facilitated by Bayesian optimisation [?], where we select the hyperparameter configuration that produces the model that maximizes  $R^2$ . The process is terminated after 1000 iterations or 2 days, whichever condition

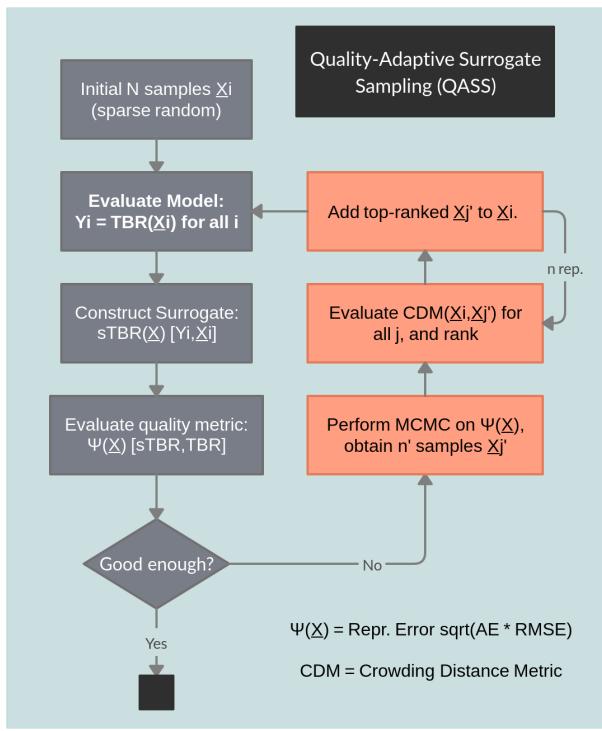
is satisfied first. The results of Experiments 1 & 2 are depicted in Figures 3 & 4 respectively, and described in Section 3.1.1.

In Experiment 3, the 20 best-performing hyperparameter configurations for each model family are used to train surrogates on sets of various sizes to investigate their scaling properties. In particular, we track the metrics from Table 3 as functions of training set size (1, 2, 5, 10, 12, 15 and 20 thousands of samples) individually for each family. This allows their comparison based on observed trends, and estimation of optimal training set sizes. The results of this experiment are shown in Figure 5 and discussed in Section 3.1.2.

Finally, Experiment 4 aims to produce surrogates suitable for practical use by retraining selected well-scaling instances on large training sets in order to satisfy the goals of this work. The results of this process are displayed in Figure 6 and in Table 5, and summarized in Section 3.1.3.

## 2.2. Adaptive Approach

Adaptive sampling techniques appear frequently in the literature and have been spe-



**Figure 2.** Schematic of QASS algorithm

cialized for surrogate modelling, where precision is implicitly limited by the quantity of training samples which are available from the expensive model. Garud's [?] “Smart Sampling Algorithm” achieved notable success by incorporating surrogate quality and crowding distance scoring to identify optimal new samples, but was only tested on a single-parameter domain. We theorized that a nondeterministic sample generation approach, built around Markov Chain Monte Carlo methods (MCMC), would fare better for high-dimensional models by more thoroughly exploring all local optima in the feature space. MCMC produces each sample point according to a jump step drawn from a shared proposal distribution. These sample points will converge to a desired posterior distribution, so long as the acceptance probability meets certain statistical criteria (see [?] for a review).

Many researchers have embedded surrogate methods into MCMC strategies for parameter optimisation [?, ?], in particular the ASMO-PODE algorithm [?] which makes use of MCMC-based adaptive sampling. Our approach draws inspiration from ASMO-PODE, but instead uses MCMC to generate samples which increase surrogate precision throughout the entire parameter space.

We designed the Quality-Adaptive Surrogate Sampling algorithm (QASS, Figure 2) to iteratively increment the training/test set with sample points which maximize surrogate error and minimize a crowding distance metric (CDM) [?]. **Error maximization is desirable for these sample points because it identifies regions of parameter space where the surrogate most needs to be improved.** On each iteration following an initial training of the surrogate on  $N$  uniformly random samples, the surrogate was trained and absolute error calculated. MCMC was then performed on the error function generated by performing nearest-neighbor interpolation on these test error points. The resultant samples were culled by 50% according to the CDM, and then the  $n$  highest-error candidates were selected for reintegration with the training/test set, beginning another training epoch. Validation was also performed during each iteration on independent, uniformly-random sample sets.

### 3. Results

#### 3.1. Decoupled Approach

**3.1.1. Hyperparameter Tuning** The results displayed in Figure 3 indicate that in the first, simplified case GBTs clearly appear to be the most accurate as well as the fastest surrogate family in terms of mean prediction time. Following that, we note that ERTs,

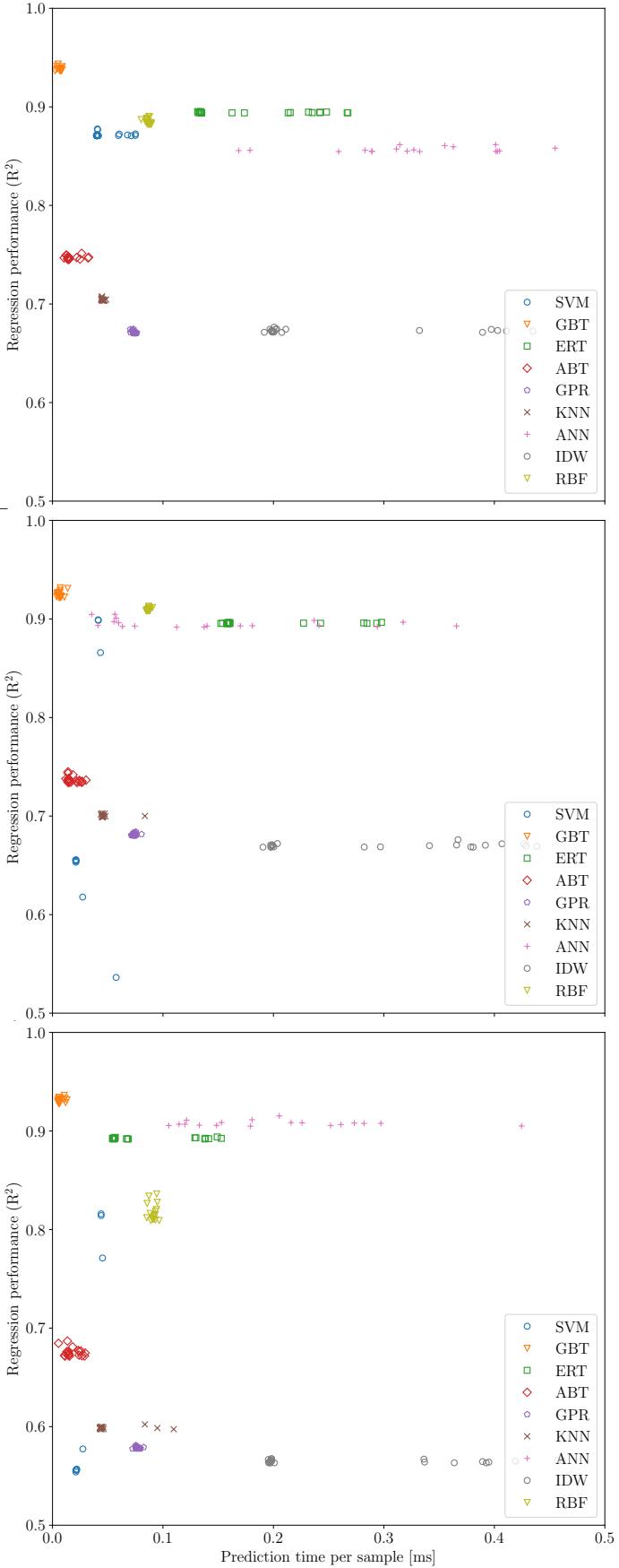
SVMs and ANNs also achieved satisfactory results with respect to both examined metric. In addition, prediction times of GBTs and SVMs show relatively lower variance than those of ERTs and ANNs. Even though the remainder of tested surrogate families does not exhibit prohibitive complexity, its regression performance falls below average.

BBM	BCM	BMM	BSM	FCM	FSM
Li <sub>4</sub> SiO <sub>4</sub>	H <sub>2</sub> O	Be <sub>12</sub> Ti	eurofer	H <sub>2</sub> O	eurofer
Li <sub>4</sub> SiO <sub>4</sub>	He	Be <sub>12</sub> Ti	eurofer	H <sub>2</sub> O	eurofer
Li <sub>4</sub> SiO <sub>4</sub>	H <sub>2</sub> O	Be <sub>12</sub> Ti	eurofer	He	eurofer
Li <sub>4</sub> SiO <sub>4</sub>	He	Be <sub>12</sub> Ti	eurofer	He	eurofer

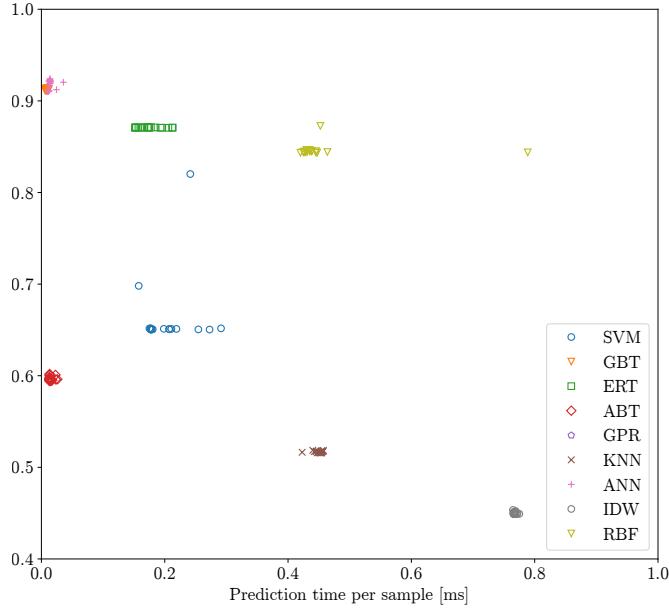
**Table 4.** Slices 1-4 of the domain space (discrete parameter assignments) explored in Experiment. Columns correspond to abbreviated parameter names listed in Table 1.

Comparing these results with those from the second, unrestricted experiment (shown in Figure 4), we observe that many surrogate families consistently underperform. The least affected models appear to be GBTs, ANN and ERTs, which are known to be capable of capturing relationships involving mixed feature types that were deliberately withheld in the first experiment. With only negligible difference, the first two of these families appear to be tied for the best performance as well as the shortest prediction time. We observe that ERTs and RBFs also demonstrated satisfactory results, clearly outperforming the remaining surrogates in terms of regression performance, and in some cases also in prediction time.

Following both hyperparameter tuning experiments, we conclude that while domain restrictions employed in the first case have proved effective in improving the regression performance of some methods, their performance fluctuates considerably depending on the selected slices. For instance, the variance in SVM performance in slice 1 is much lower than in slices



**Figure 3.** Experiment 1 results. 20 best surrogates per each considered family, plotted in terms of  $\bar{t}_{\text{pred}}$  and  $R^2$  with 3 selected slices out of 4 (defined in Table 4).

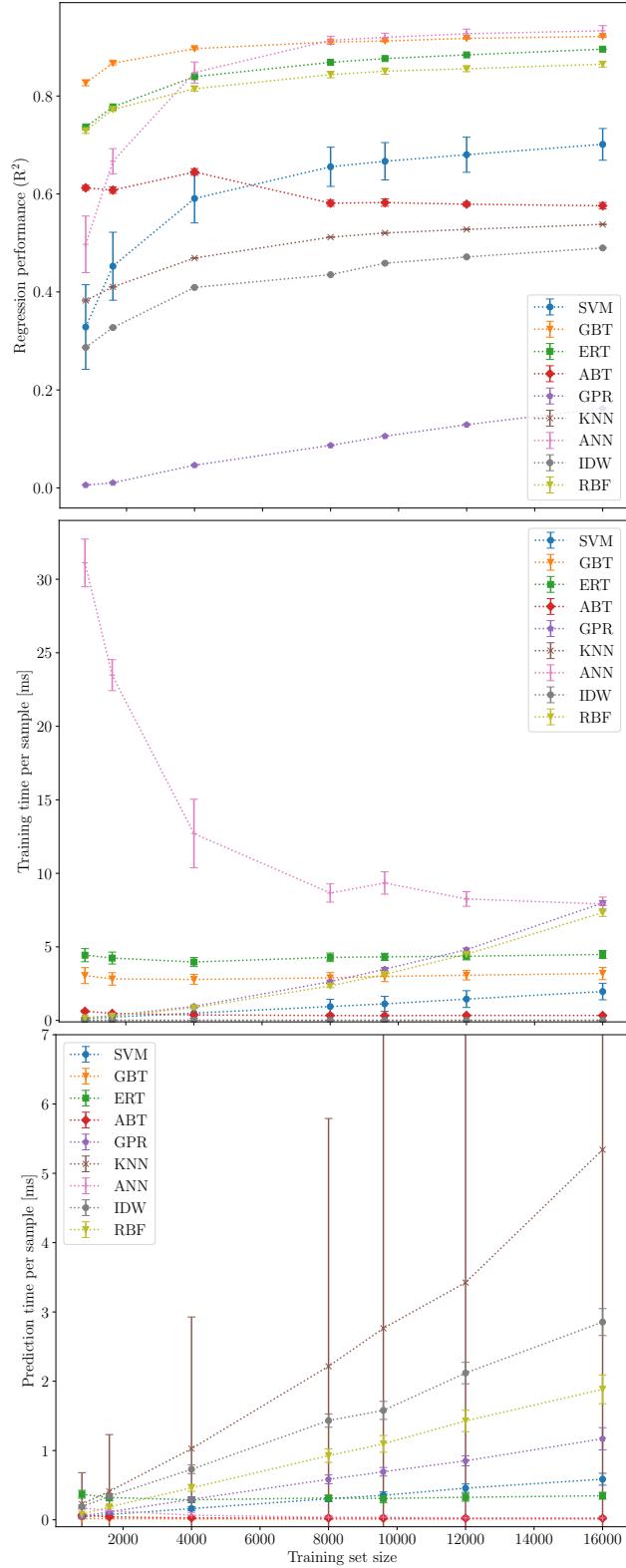


**Figure 4.** Experiment 2 results, plotted analogously to Figure 3.

2-3, and both KNNs and ABTs perform much better in slices 1-2 than in slice 3. Furthermore, in all instances the best results are achieved by families of surrogates that do not benefit from this restriction (GBTs, ANNs, ERTs).

**3.1.2. Scaling Benchmark** The results shown in Figure 5 suggest that in terms of regression performance the most accurate families from the previous experiments consistently maintain their relative advantage over others, even as more training points are introduced. While such families achieve nearly comparable performance on the largest dataset, in the opposite case tree-based ensemble approaches clearly outperform ANNs. For instance, GBTs achieve  $MAE = 0.107$  that is nearly half of the  $MAE = 0.186$  reached by ANNs, representing a clear benefit given vastly disparate training and prediction times. This trend continues to set sizes up to 6000.

Consistent with our expectations, the shortest training times were achieved by



**Figure 5.** Experiment 3 results, displayed as a function of  $N_0$ . From top to bottom,  $R^2$ ,  $\bar{t}_{\text{trn.}}$ ,  $\bar{t}_{\text{pred.}}$ .

instance-based learning methods (KNN, IDW) that are trained trivially at the expense of increased lookup complexity later during prediction. Furthermore, we observe that the majority of tree-based ensemble algorithms also perform and scale well, unlike RBFs and GPR which appear to behave superlinearly. We note that ANNs, which are the only family to utilize parallelisation during training, show an inverse scaling characteristic. We suspect that this effect may be caused by a constant multi-threading overhead that dominates the training process on relatively small sets.

Finally, all tested families with the exception of previously mentioned instance-based models offer desirable prediction times. Analogous to previous experiments, GBTs, ABTs and ANNs appear to be tied, as they not only exhibit comparable times but also similar scaling slopes. Following those, we note a clear hierarchy of ERTs, SVMs, GPR and RBFs, trailed by IDW and KNNs.

*3.1.3. Model Comparison* In Experiment 4, we aim to create models that yield: (a) the best regression performance regardless of other features, (b) acceptable performance with the shortest mean prediction time, or (c) acceptable performance with the smallest training set. To this end, we trained 8 surrogates that are presented in Figure 6 and Table 5. We compared these surrogates with the baseline represented by Paramak per-sample evaluation time  $\bar{t}_{\text{eval.}} = 7.777 \pm 2.810$  s, which was measured earlier on a set of 500000 samples.

Having selected ANNs, GBTs, ERTs, RBFs and SVMs based on the results of Experiments 2 & 3, we utilized the best-performing hyperparameters. In pursuit of goal (a), the best approximator (Model 1, ANN) achieved  $R^2 = 0.998$  and mean prediction time  $\bar{t}_{\text{pred.}} = 1.124$   $\mu$ s. These correspond to

a standard error  $S = 0.013$  and a relative speedup  $\omega = 6.92 \times 10^6$  with respect to Paramak. Satisfying goal (b), the fastest model (Model 2, ANN) achieved  $R^2 = 0.985$ ,  $\bar{t}_{\text{pred.}} = 0.898$   $\mu$ s,  $S = 0.033$  and  $\omega = 8.66 \times 10^6$ . While these surrogates were trained on the entire available set of 500000 datapoints, to satisfy goal (c) we also trained a more simplified model (Model 4, GBT) that achieved  $R^2 = 0.913$ ,  $\bar{t}_{\text{pred.}} = 6.125$   $\mu$ s,  $S = 0.072$  and  $\omega = 1.27 \times 10^6$  with a set of size only 10000.

Overall we found that due to their superior performance, boosted tree-based approaches seem to be advantageous for fast surrogate modelling on relatively small training sets (up to the order of  $10^4$ ). Conversely, while neural networks perform poorly in such a setting, they dominate on larger training sets (at least of the order of  $10^5$ ) both in terms of regression performance and mean prediction time.

### 3.2. Adaptive Approach

In order to test our QASS prototype, several functional toy theories for TBR were developed as alternatives to the expensive MC model.

QASS performance was verified by training an ANN on these theories for varied quantities of initial, incremental, and MCMC candidate samples. By far the most useful of these was the following sinusoidal theory, as ANNs trained on this model demonstrated similar performance to those on the expensive MC model:

$$\text{TBR} = |C|^{-1} \sum_{i \in C} [1 + \sin(2\pi n(x_i - 1/2))] \quad (3)$$

where  $C$  denotes the continuous parameter space, and  $n$  is an adjustable wavenumber parameter.

An increase in initial samples with increment held constant had a strong impact on final surrogate precision, an early confirmation of basic functionality. An increase in MCMC

**Table 5.** Results of Experiment 4. Here, means and standard deviations are reported over 5 cross-validation folds,  $|\mathcal{T}|$  denotes cross-validation set size ( $\times 10^3$ ) and  $\omega$  is a relative speedup with respect to  $\bar{t}_{\text{eval.}} = 7.777 \pm 2.810$  s measured in Paramak over 500000 samples. The best-performing method(s) under each metric are highlighted in bold.

Model	$ \mathcal{T} $	Regression performance			Computational complexity			
		MAE [TBR]	$S$ [TBR]	$R^2$ [rel.]	$R^2_{\text{adj.}}$ [rel.]	$\bar{t}_{\text{trn.}}$ [ms]	$\bar{t}_{\text{pred.}}$ [ms]	$\omega$ [rel.]
1 (ANN)	500	<b><math>0.009 \pm 0.000</math></b>	<b><math>0.013 \pm 0.001</math></b>	<b><math>0.998 \pm 0.000</math></b>	<b><math>0.998 \pm 0.000</math></b>	$3.659 \pm 0.035$	<b><math>0.001 \pm 0.000</math></b>	$6916416 \times$
2 (ANN)	500	$0.025 \pm 0.001$	$0.033 \pm 0.001$	$0.985 \pm 0.001$	$0.985 \pm 0.001$	$2.989 \pm 0.026$	<b><math>0.001 \pm 0.000</math></b>	<b><math>8659251 \times</math></b>
3 (GBT)	200	$0.058 \pm 0.001$	$0.059 \pm 0.000$	$0.941 \pm 0.001$	$0.941 \pm 0.001$	$2.221 \pm 0.010$	$0.007 \pm 0.000$	$1169933 \times$
4 (GBT)	10	$0.071 \pm 0.002$	$0.072 \pm 0.003$	$0.913 \pm 0.006$	$0.912 \pm 0.006$	<b><math>1.621 \pm 0.008</math></b>	$0.006 \pm 0.000$	$1269777 \times$
5 (ERT)	200	$0.051 \pm 0.000$	$0.056 \pm 0.000$	$0.950 \pm 0.001$	$0.950 \pm 0.001$	$2.634 \pm 0.010$	$0.214 \pm 0.004$	$36308 \times$
6 (ERT)	40	$0.068 \pm 0.000$	$0.072 \pm 0.000$	$0.917 \pm 0.001$	$0.917 \pm 0.001$	$2.368 \pm 0.005$	$0.188 \pm 0.008$	$41370 \times$
7 (RBF)	50	$0.068 \pm 0.001$	$0.077 \pm 0.002$	$0.910 \pm 0.003$	$0.910 \pm 0.003$	$3.453 \pm 0.019$	$1.512 \pm 0.016$	$5143 \times$
8 (SVM)	200	$0.062 \pm 0.000$	$0.094 \pm 0.002$	$0.891 \pm 0.003$	$0.891 \pm 0.003$	$33.347 \pm 0.382$	$2.415 \pm 0.011$	$3220 \times$

candidate samples was seen to have a positive but very weak effect on final surrogate precision, suggesting that the runtime of MCMC on each iteration could be limited for increased efficiency. We also found that an optimum increment exists and is monotonic with initial sample quantity, above or below which models showed slower improvement on both the training and evaluation sets, and a larger minimum error on the evaluation set. This performance distinction will be far more significant for an expensive model such as Paramak, where the number of sample evaluations is the primary computational bottleneck.

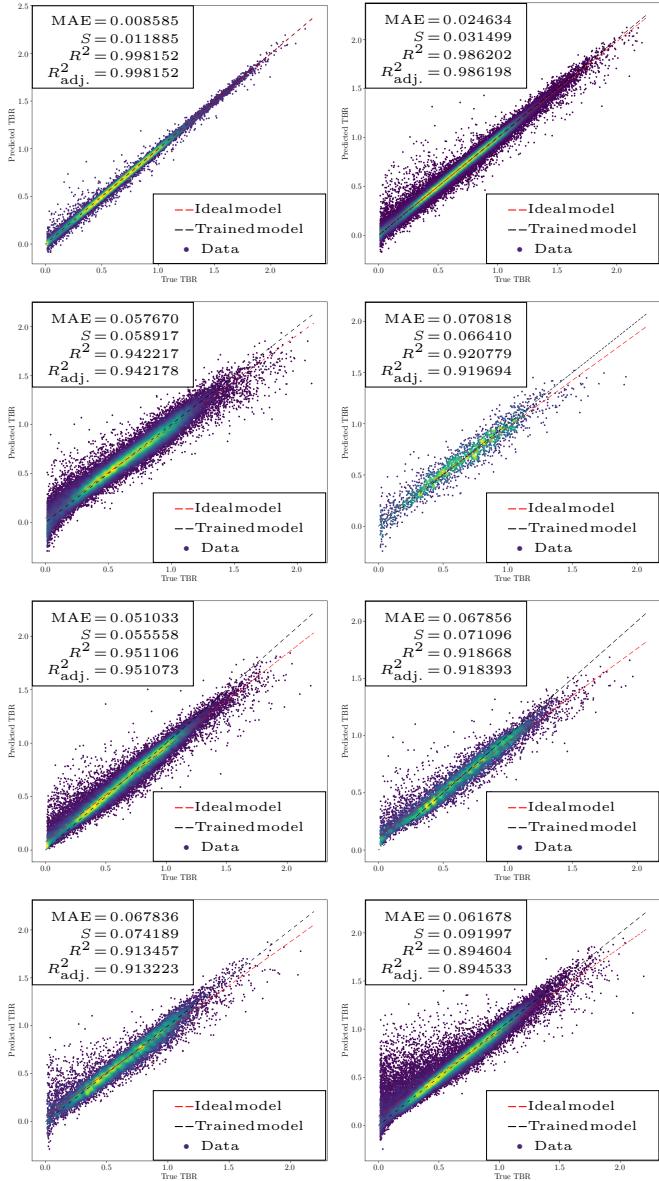
A plateau effect in surrogate error on the evaluation set was universal to all configurations, and initially suspected to be a residual effect of retraining the same ANN instance without adjustment to data normalisation. A “Goldilocks scheme” for checking normalisation drift was implemented and tested, but did not affect QASS performance. Schemes in which the ANN is periodically retrained were also discarded, as the retention of network weights from one iteration to the next was demonstrated to greatly benefit QASS efficiency. Further insight came from direct comparison between QASS and a baseline scheme with uniformly random incremental samples, shown in Figure 7.

Such tests revealed that while QASS has

unmatched performance on its own adaptively-sampled training set, it is outperformed by the baseline scheme on uniformly-random evaluation sets. We inferred that while QASS excels in learning the most strongly peaked regions of the TBR theory, this comes at the expense of precision in broader, smoother regions where uniformly random sampling suffices. Therefore a mixed scheme was implemented, with half MCMC samples and half uniformly-random samples incremented on each iteration, which is also shown in Figure 7. An increase in initial sample size was also observed to improve precision in these smooth regions of the toy theory, as the initial samples were uniformly-random. As shown in Figure 8, with 100000 initial samples it was possible to obtain a  $\sim 40\%$  decrease in error as compared to the baseline scheme, from 0.0025 to 0.0015 mean averaged error. Comparing at the point of termination for QASS, this corresponds to a  $\sim 6\%$  decrease in the number of total samples needed to train a surrogate model while achieving the same error.

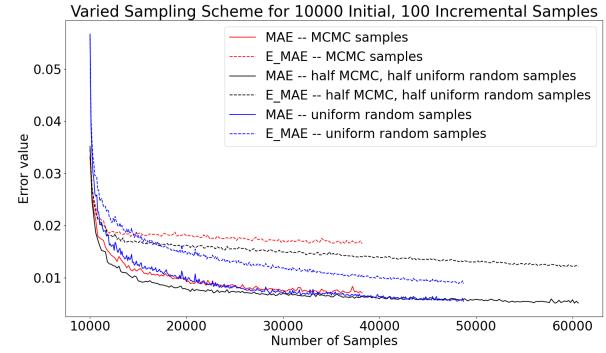
#### 4. Conclusion

We employed a broad spectrum of data analysis and machine learning techniques to develop a library of fast and high-quality surrogate

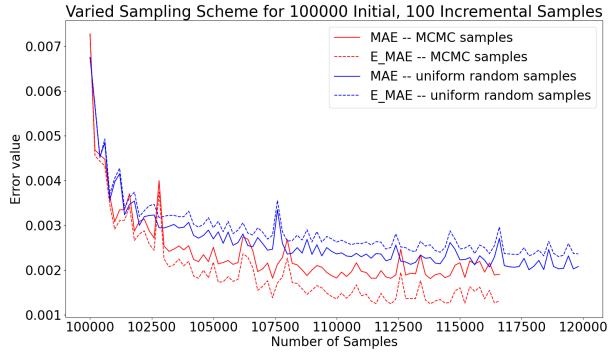


**Figure 6.** Regression performance of Models 1-8 (from left to right, top to bottom) in Experiment 4, viewed as true vs. predicted TBR on a test set of a selected cross-validation fold. Points are colored by density.

models for the expensive Paramak TBR model. After reviewing 9 surrogate model families, examining their behaviour on a constrained and unrestricted feature space, and studying their scaling properties, we retrained the best-performing instances to produce properties desirable for practical use. The fastest surrogate, an artificial neural network trained



**Figure 7.** Absolute training error for QASS, baseline scheme, and mixed scheme.



**Figure 8.** Absolute training error for QASS and baseline scheme, with 100000 initial samples.

on 500000 datapoints, attained an  $R^2 = 0.985$  with mean prediction time of  $0.898 \mu\text{s}$ , representing a relative speedup of  $8 \cdot 10^6$  with respect to Paramak. Furthermore, we demonstrated the possibility of achieving comparable results using only a training set of size 10000.

We additionally developed a novel adaptive sampling algorithm, QASS, capable of interfacing with any of the [surrogate models presented in this work](#). Preliminary testing on a toy theory, qualitatively comparable to Paramak, demonstrated the effectiveness of QASS and [key](#) behavioral trends. With 100000 initial samples and 100 incremental samples per iteration, QASS achieved a  $\sim 40\%$  decrease in surro-

gate error compared to a baseline random sampling scheme. Further optimisation over the hyperparameter space has strong potential to increase this performance by further reduction of necessary expensive samples, in particular by decreasing the required quantity of initial samples. This will allow for future deployment of QASS on top of any of **our most effective identified TBR surrogate models**.

Relevant source code, model instances and datasets are freely available online as well as a more detailed technical report [?, ?].

## 5. Acknowledgements

PM & GVG were supported the STFC UCL Centre for Doctoral Training in Data Intensive Science (grant no. ST/P006736/1).

This project was supported by the EU Horizon 2020 research & innovation programme [grant No 758892, ExoAI]. N. Nikolaou acknowledges the support of the NVIDIA Corporation’s GPU grant.

This work has been carried out within the framework of the EUROfusion consortium and has received funding from the Euratom research and training programme 2014-2018 and 2019-2020 under grant agreement No 633053. The views and opinions expressed herein do not necessarily reflect those of the European Commission.

This work has been partly funded by the Institutional support for the development of a research organization (DKRVO, Czech Republic).

This work has also been part-funded by the RCUK Energy Programme (grant number EP/I501045/1).

## 6. References