# Multi-Stage Multi-Fidelity Gaussian Process Modeling, with Application to Heavy-Ion Collisions

Yi Ji[1,5],  Henry Shaowu Yuchi[2,5]

Derek Soeder[3],  J.-F. Paquet[3,4]

Steffen A. Bass[3],  V. Roshan Joseph[2],  C. F. Jeff Wu[2],  Simon Mak[1,*]

September 29, 2022

## Abstract

In an era where scientific experimentation is often costly, multi-fidelity emulation provides a powerful tool for predictive scientific computing. While there has been notable work on multi-fidelity modeling, existing models do not incorporate an important multi-stage property of multi-fidelity simulators, where multiple fidelity parameters control for accuracy at different experimental stages. Such multi-stage simulators are widely encountered in complex nuclear physics and astrophysics problems. We thus propose a new Multi-stage Multi-fidelity Gaussian Process ($M^2GP$) model, which embeds this multi-stage structure within a novel non-stationary covariance function. We show that the $M^2GP$ model can capture prior knowledge on the numerical convergence of multi-stage simulators, which allows for cost-efficient emulation of multi-fidelity systems. We demonstrate the improved predictive performance of the $M^2GP$ model over state-of-the-art methods in a suite of numerical experiments and two applications, the first for emulation of cantilever beam deflection and the second for emulating the evolution of the quark-gluon plasma, which was theorized to have filled the Universe shortly after the Big Bang.

*Keywords: Bayesian Nonparametrics, Multi-Fidelity Emulation, Multi-Stage Simulation, Quark-Gluon Plasma, Surrogate Modeling.*

---

[1]Department of Statistical Science, Duke University

[2]H. Milton Stewart School of Industrial & Systems Engineering, Georgia Institute of Technology

[3]Department of Physics, Duke University

[4]Department of Physics and Astronomy & Department of Mathematics, Vanderbilt University

[5]Joint first authors

[*]Corresponding author

# 1 Introduction

Computer experimentation is widely used for modeling complex scientific and engineering systems, particularly when physical experiments are costly, unethical, or impossible to perform. This shift from physical to computer experimentation has found success in a wide range of physical science applications, from rocket design (Mak et al. 2018), solar irradiance modeling (Sun et al. 2019) to 3D printing (Chen et al. 2021). However, as systems become more complex and realistic, such computer experiments also become more expensive, thus placing a heavy computational burden on design exploration and optimization. Statistical *emulators* (Santner et al. 2003) have shown great promise in tackling this limitation. The idea is simple but effective: computer experiments are first performed at carefully chosen design points, then used as training data to fit an *emulator* model to efficiently predict and quantify uncertainty on the expensive virtual experiment.

In recent years, however, with the increasing sophistication of modern scientific problems, an emerging challenge for emulators is the simulation of high-fidelity training data, which can be prohibitively expensive. One way to address this is via *multi-fidelity emulation*, which makes use of training simulation data of multiple *fidelities* (or accuracy) for model fitting. Such multi-fidelity data can often be generated by varying different *fidelity parameters*, which control the precision of the numerical experiment. There are a wide variety of fidelity parameters, ranging from mesh sizes for finite element analysis (Park et al. 1997, More & Bindu 2015) to time-steps for dynamical system simulation (Vanden-Eijnden 2003). The goal is to leverage information from lower-fidelity (but cheaper) simulations to enhance predictions for the high-fidelity (but expensive) model, thus allowing for improved emulation and uncertainty quantification at lower computational costs.

There has been much recent work on multi-fidelity emulation, particularly for Gaussian process (GP) modeling. This includes the seminal work of Kennedy & O'Hagan (2000), which presented a first-order autoregressive model for integrating information over a hierarchy of simulation models, from lowest to highest fidelity. This so-called Kennedy-O'Hagan model has then been extended in various works, including a Bayesian hierarchical im-
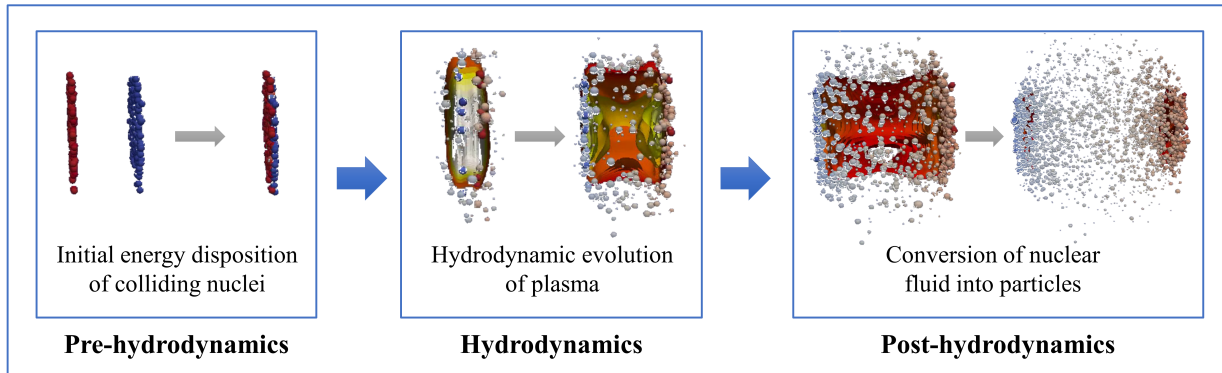
Figure 1: Three-stage simulation of the quark-gluon plasma.

plementation in Qian et al. (2006), the multi-fidelity optimization approach in Forrester et al. (2007), and the nonlinear fusion model in Perdikaris et al. (2017). Tuo et al. (2014) proposed a multi-fidelity emulator for finite element analysis (FEA), which considers the discretization mesh size as a single fidelity parameter. Such models have been widely applied in engineering design and scientific computing; see, e.g., Kou & Zhang (2019), Shi et al. (2020), Jin et al. (2021), Liyanage et al. (2022). Furthermore, experimental design for multi-fidelity emulators has been explored (Xiong et al. 2013), including a Bayesian optimization sequential design strategy in He et al. (2017).

The above methods, however, have limitations when applied to our motivating nuclear physics application, which we describe below. Here, we are interested in the Quark-Gluon Plasma (QGP), a deconfined phase of nuclear matter consisting of elementary quarks and gluons. The QGP was theorized to have filled the Universe shortly after the Big Bang, and the study of this plasma sheds light on the properties of this unique phase of matter. This plasma can be simulated at a small scale by virtually colliding heavy ions together at near-light speeds in particle colliders. Simulating such collisions requires a *multi-stage* system of complex dynamical models to faithfully capture the detailed evolution of the plasma. Consider in particular the three-stage simulation framework in Everett et al. (2021) (see also (Gale et al. 2013, Heinz & Snellings 2013, De Souza et al. 2016)), which models the initial energy disposition of the heavy ions, the hydrodynamic evolution of the plasma

3

after the collision, and the subsequent conversion of nuclear fluid into particles. Figure 1 visualizes this multi-stage procedure. At each stage, the simulation of the component physics can involve *multiple* and *different* fidelity parameters, controlling, e.g., the size of the hydrodynamics' spatial mesh, or the time-scale for dynamic evolution.

This *multi-stage multi-fidelity* framework, which is widely encountered in complex physics simulators, poses several challenges for existing multi-fidelity emulators. First, since there are multiple fidelity parameters to set for each simulation stage, the resulting simulation runs typically cannot be ranked from lowest to highest fidelity, which is required for a direct application of Kennedy-O'Hagan-type models. For example, to gauge the effects of three fidelity parameters, the physicist may choose to run the simulator in three different ways, each with higher fidelity at one stage and lower fidelity at remaining stages. A priori, it is unclear if these three simulation approaches can be ranked from lowest to highest fidelity. Second, unlike the multi-fidelity emulator in Tuo et al. (2014) (which allows only one fidelity parameter), there are *multiple* fidelity parameters which should be accounted for when training emulators with multi-stage simulations. Neglecting this multi-stage structure for emulation can result in significantly poorer predictive performance, as we show later. A broader emulation model is thus needed to tackle the challenges presented by this multi-stage multi-fidelity framework, which is widely encountered in nuclear physics (Ji et al. 2021) and astrophysics (Ho et al. 2022).

We propose in this work a new GP emulator which addresses these challenges. The proposed *Multi-stage Multi-fidelity Gaussian Process* ($M^2$GP) model makes use of a novel *non-stationary* covariance function, which captures prior information on the numerical convergence of multi-stage simulators. By embedding this prior knowledge within the emulator, we show that the $M^2$GP can indeed yield improved emulation performance and uncertainty quantification over existing methods, in a suite of numerical experiments, a beam deflection problem in finite element analysis, and an application to the motivating heavy-ion collision problem. Section 2 reviews several existing multi-fidelity emulators and outlines the motivating QGP problem. Section 3 presents the model specification for the proposed

4

M²GP emulator. Section 4 discusses implementation details for the M²GP, including parameter estimation and experimental design. Section 5 compares the proposed model with existing methods on a suite of numerical experiments. Finally, Section 6 demonstrates the effectiveness of the M²GP for the motivating QGP application as well as a cantilever beam deflection problem. Section 7 concludes the paper.

# 2 Preliminaries & Motivation

In this section, we first provide a brief overview of the Gaussian process model. Next, we review the Kennedy-O'Hagan model and the multi-fidelity model in Tuo et al. (2014), then discuss the limitations of such models for the aforementioned QGP application, thus motivating the proposed M²GP modeling framework.

## 2.1 Gaussian Process

Gaussian process (GP) modeling is a popular Bayesian nonparametric approach for supervised learning (Williams & Rasmussen 2006), with broad applications for computer experiments (Santner et al. 2003). The specification of a GP model involves two key ingredients: the mean function and the covariance function. Let $\mathbf{x} \in [0,1]^p$ be the input parameters (sufficiently scaled) for the simulator, and let $\eta(\mathbf{x})$ be the corresponding output of the simulator. A GP model places the following prior on the unknown response surface $\eta(\cdot)$:

$$\eta(\cdot) \sim \mathcal{GP}(\mu(\cdot), k(\cdot, \cdot)). \tag{1}$$

Here, $\mu(\cdot)$ is the mean function controlling the centrality of the stochastic process; this is typically set to be a constant in the absence of prior knowledge on mean trends. The function $k(\cdot, \cdot)$ is the covariance function that controls the smoothness of its sample paths. Common choices of $k(\cdot, \cdot)$ include the squared-exponential and Matérn kernels (Santner et al. 2003).

Let $\mathcal{D} = \{\mathbf{x}_1, \cdots, \mathbf{x}_n\}$ denote the simulated input points, and $\mathbf{y} = [\eta(\mathbf{x}_1), \cdots, \eta(\mathbf{x}_n)]$

be the simulated outputs. Assuming that the kernel hyperparameters are fixed and known (we will discuss the estimation of such parameters later in Section 4.1), the predictive distribution $\eta(\mathbf{x}^*)$ at the new input $\mathbf{x}^*$ conditional on data $\{\mathcal{D}, \mathbf{y}\}$ is given by:

$$\eta(\mathbf{x}^*)|\mathcal{D}, \mathbf{y} \sim \mathcal{GP}(\hat{\mu}(\mathbf{x}^*), s^2(\mathbf{x}^*)). \tag{2}$$

Here, the posterior mean and variance are given by:

$$\begin{aligned}
\hat{\mu}(\mathbf{x}^*) &= \mu(\mathbf{x}^*) + \mathbf{k}(\mathbf{x}^*, \mathcal{D})^T \mathbf{K}(\mathcal{D})^{-1}(\mathbf{y} - \boldsymbol{\mu}(\mathcal{D})), \\
s^2(\mathbf{x}^*) &= k(\mathbf{x}^*, \mathbf{x}^*) - \mathbf{k}(\mathbf{x}^*, \mathcal{D})^T \mathbf{K}(\mathcal{D})^{-1}\mathbf{k}(\mathbf{x}^*, \mathcal{D}),
\end{aligned} \tag{3}$$

where $\mathbf{k}(\mathbf{x}^*, \mathcal{D}) = [k(\mathbf{x}^*, \mathbf{x}_1), \cdots, k(\mathbf{x}^*, \mathbf{x}_n)]$ is the vector of covariances, $\boldsymbol{\mu}(\mathcal{D}) = [\mu(\mathbf{x}_1), \cdots, \mu(\mathbf{x}_n)]$ is the vector of means, and $\mathbf{K}(\mathcal{D})$ is the covariance matrix for the training data. The models introduced later in this paper will make use of these closed-form predictive equations with different choices of covariance functions.

## 2.2  The Kennedy-O'Hagan model

In the seminal work of Kennedy & O'Hagan (2000), the authors proposed a first-order autoregressive model for linking outputs from a hierarchy of $K$ simulators, from lowest fidelity (simulator 1) to highest fidelity (simulator $K$). Let $\eta_k(\mathbf{x})$ denote the output from simulator $k$ at standardized input parameters $\mathbf{x} \in [0, 1]^p$. The Kennedy-O'Hagan (KOH) model is specified as:

$$\eta_k(\mathbf{x}) = \rho_{k-1}\eta_{k-1}(\mathbf{x}) + \delta_k(\mathbf{x}), \quad k = 2, \cdots, K. \tag{4}$$

Here, $\rho_{k-1}$ is a regression scale factor, and $\delta_k(\mathbf{x})$ is a bias term which models for the discrepancy between simulator $k-1$ and $k$. The bias term $\delta_k(\mathbf{x})$ can then modeled by a stationary GP with a squared-exponential covariance function (Santner et al. 2003) :

$$\text{Cov}\left[\delta_k(\mathbf{x}), \delta_k(\mathbf{x}')\right] = \sigma_k^2 \exp\left\{-\sum_{i=1}^{p} \phi_{k,i}(\mathbf{x}_i - \mathbf{x}'_i)^2\right\}, \tag{5}$$

where $\phi_{k,i}$ is the weight parameter for the $i^{\text{th}}$ input parameter at the $k^{\text{th}}$ fidelity level. Such a model allows one to integrate information from a sequence of simulator models with varying fidelity levels, to efficiently emulate the highest-fidelity simulator model.

The KOH multi-fidelity model has subsequently been extended in a variety of ways, including a Bayesian implementation in Qian & Wu (2008) and a nonlinear extension in Perdikaris et al. (2015); see also Reese et al. (2004), DiazDelaO & Adhikari (2012), Fricker et al. (2013). This modeling framework is also closely related to the idea of co-kriging (Stein & Corsten 1991), which is widely used in spatial statistics. However, the aforementioned methods assume that the multi-fidelity training data can be *ranked* from lowest to highest fidelity. As such, this body of literature does not directly apply to the motivating problem of multi-stage multi-fidelity emulation, where simulation accuracy is controlled by *multiple* fidelity parameters, and thus there is no clear ranking of training data from lowest to highest fidelity. There are several ways to force existing models on this problem, but each has its shortcomings. One could design the data such that the training simulations are ranked (e.g., increasing all fidelity parameters simultaneously), but this would result in highly inefficient designs which fail to sufficiently explore the space of fidelity parameters. One could also arbitrarily assign a *single* "artificial" fidelity level for each simulation, which imposes a ranking on the training runs. This, however, *neglects* the rich multi-stage multi-fidelity framework (i.e., the "science") for the simulator, which can lead to significantly poorer predictive performance from the emulator, as we show later.

## 2.3   The Tuo-Wu-Yu model

For problems where the fidelity level is controlled by a *single* continuous fidelity parameter $t$ (e.g., mesh size), Tuo et al. (2014) proposed an alternate model (we call this the TWY model) which can make use of such information. Let $\eta(\mathbf{x}, t)$ denote the deterministic code output at standardized inputs $\mathbf{x} \in [0, 1]^p$ and at fidelity parameter $t$. Here, $t$ is typically assumed to be between 0 and 1, with a smaller $t$ indicating a finer mesh size or, equivalently,

higher mesh density. The TWY model adopts the following model for $\eta(\mathbf{x}, t)$:

$$\eta(\mathbf{x}, t) = \eta(\mathbf{x}, 0) + \delta(\mathbf{x}, t) =: \phi(\mathbf{x}) + \delta(\mathbf{x}, t). \tag{6}$$

Here, $\phi(\mathbf{x}) := \eta(\mathbf{x}, 0)$ denotes the "exact" simulation output at input $\mathbf{x}$ at the highest (limiting) fidelity $t = 0$, and $\delta(\mathbf{x}, t)$ denotes the discrepancy (or bias) between this exact solution and realized simulation output with mesh size $t$. In practical problems, the exact solution $\eta(\mathbf{x}, 0)$ is typically *not obtainable* numerically, since some level of approximation (e.g., mesh or time discretization) is needed for simulating the system. The goal is to leverage simulation training data of the form $\{\eta(\mathbf{x}_i, t_i)\}_{i=1}^n$ along with an appropriate model on (6) to predict the exact solution $\eta(\mathbf{x}, 0)$.

Since $\phi(\mathbf{x})$ and $\delta(\mathbf{x}, t)$ are unknown *a priori*, these terms are modeled in Tuo et al. (2014) by two independent Gaussian processes. For $\phi(\mathbf{x})$, a standard GP prior is assigned with constant mean and a stationary correlation (e.g., squared-exponential) function. For the bias term $\delta(\mathbf{x}, t)$, a *non-stationary* zero-mean GP prior is assigned, with covariance function:

$$\text{Cov}[\delta(\mathbf{x_1}, t_1), \delta(\mathbf{x_2}, t_2)] = \sigma_2^2 K_{\mathbf{x}}^{\delta}(\mathbf{x_1}, \mathbf{x_2}) \min(t_1, t_2)^l, \tag{7}$$

where $K_{\mathbf{x}}^{\delta}(\cdot, \cdot)$ is a stationary correlation function on input parameters $\mathbf{x}$. One can view this as a product of two kernels, where the kernel on the fidelity parameter $t$ is non-stationary and closely resembles that of a Brownian motion.

This choice of non-stationary kernel over the single fidelity parameter $t$ can be reasoned from a Bayesian modeling perspective. Consider the GP model with covariance function (7) as a prior model on discrepancy $\delta(\mathbf{x}, t)$. Prior to data, one can show that:

$$\lim_{t \to 0} \delta(\mathbf{x}, t) = 0, \quad \text{for all } \mathbf{x} \in [0, 1]^p \text{ almost surely.} \tag{8}$$

The TWY model thus assumes a priori that the discrepancy term should converge to 0 as fidelity parameter $t$ goes to 0, or equivalently, the simulation output $\eta(\mathbf{x}, t)$ converges to the exact solution $\phi(\mathbf{x})$, as we increase the fidelity of the simulator. This can be seen as a way of integrating *prior* information on the numerical convergence of the simulator within

the prior specification of the emulator model. One can further set the kernel parameter $l$ to capture additional information on known numerical convergence rates of the simulator; see Tuo et al. (2014) for details.

For the target problem of *multi-stage* multi-fidelity emulation, however, where *multiple* fidelity parameters are present, the TWY model needs to be further extended. A simple modification might be to first assign for each simulation run an "artificial" fidelity, e.g., the average of the multiple fidelity parameters, then use this single aggregate fidelity level with the TWY model for multi-fidelity emulation. However, such an approach ignores the rich multi-stage structure of the simulation framework, which can lead to poor predictive performance. We show later that, by integrating directly the multi-stage multi-fidelity nature of the simulation framework (i.e., the "science") within the M²GP, we can achieve significantly improved predictive performance in numerical experiments and for the motivating nuclear physics application.

# 3  The M²GP model

Given these limitations, we now present the proposed M²GP model for efficient emulation of multi-stage multi-fidelity simulations. Our model adopts a novel non-stationary Gaussian process model which captures *prior* information on the numerical convergence behavior of *multi-stage* simulators. Below, we outline the general M²GP model specification, then present two choices of non-stationary covariance functions which capture this desired prior information.

Let $\mathbf{x} \in [0, 1]^p$ be the vector of $p$ standardized simulation inputs for the computer code (again assumed to be deterministic), and suppose there are $k$ fidelity parameters (denoted by $\mathbf{t} \in [0, 1]^k$) which control simulation accuracy in the code. These may, e.g., consist of different mesh sizes for domain discretization and time steps at different simulation stages. As before, a smaller fidelity parameter $t_r$ (with other fidelity parameters held constant) yields more accurate simulations at higher computational costs, with $t_r = 0$ denoting the highest (limiting) fidelity level. Let $\eta(\mathbf{x}, \mathbf{t})$ denote the deterministic code output at inputs

9

$\mathbf{x}$ and fidelity parameters $\mathbf{t}$. The M$^2$GP model assumes the following decomposition of $\eta(\mathbf{x}, \mathbf{t})$:

$$\eta(\mathbf{x}, \mathbf{t}) = \eta(\mathbf{x}, \mathbf{0}) + \delta(\mathbf{x}, \mathbf{t}) := \phi(\mathbf{x}) + \delta(\mathbf{x}, \mathbf{t}). \tag{9}$$

Similar to before, $\phi(\mathbf{x}) := \eta(\mathbf{x}, \mathbf{0})$ models the "exact" simulation solution at the highest (limiting) fidelity setting of $\mathbf{t} \to \mathbf{0}$, and $\delta(\mathbf{x}, \mathbf{t})$ models the numerical discrepancy (or error) between the exact solution $\phi(\mathbf{x})$ and the simulated output $\eta(\mathbf{x}, \mathbf{t})$. Since both $\phi(\mathbf{x})$ and $\delta(\mathbf{x}, \mathbf{t})$ are unknown, we again place independent Gaussian process priors on both terms. For $\phi(\mathbf{x})$, a standard GP is assigned with user-defined basis functions for the mean and a stationary correlation function. In a later implementation, we will make use of linear basis functions along with the popular squared-exponential correlation function:

$$\mathrm{Cov}[\phi(\mathbf{x}_1), \phi(\mathbf{x}_2)] = \sigma_1^2 K_{\mathbf{x}}^{\phi}(\mathbf{x}_1, \mathbf{x}_2) = \sigma_1^2 \exp\left\{ -\sum_{s=1}^{p} \gamma_s (x_{1,s} - x_{2,s})^2 \right\}, \tag{10}$$

where $\gamma_s$ is the weight parameter for the $s^{\text{th}}$ input dimension.

For the bias term $\delta(\mathbf{x}, \mathbf{t})$, we will carefully specify a new non-stationary covariance function which captures one's *prior* knowledge on the numerical convergence behavior. One desirable property of $\delta(\mathbf{x}, \mathbf{t})$ is the limiting constraint:

$$\lim_{\mathbf{t} \to \mathbf{0}} \delta(\mathbf{x}, \mathbf{t}) = 0, \quad \text{for all } \mathbf{x} \in [0, 1]^p \text{ almost surely.} \tag{11}$$

In words, for any inputs $\mathbf{x}$, the simulation output $\eta(\mathbf{x}, \mathbf{t})$ should converge to the underlying exact solution $\phi(\mathbf{x})$ when *all* fidelity parameters converge to zero, i.e., all fidelity levels are set to their highest (limiting) setting. Property (11) should thus be satisfied almost surely if the simulator enjoys theoretical convergence guarantees (e.g., weak convergence of PDE solutions) or is trusted to converge empirically. Another desirable property is that, for a fidelity parameter $t_r$ and fixed levels of the remaining fidelity parameters $\mathbf{t}_{-r} \neq \mathbf{0}$, we have:

$$\lim_{t_r \to 0} \delta(\mathbf{x}, \mathbf{t}) \neq 0, \quad \text{for all } \mathbf{x} \in [0, 1]^p \text{ almost surely.} \tag{12}$$

In words, for any inputs $\mathbf{x}$ and any positive fidelity parameters $\mathbf{t}_{-r}$, there should be non-negligible discrepancy between the simulation output $\eta(\mathbf{x}, \mathbf{t})$ and the underlying true solution $\phi(\mathbf{x})$. This is again intuitive when the variables in $\mathbf{t}_{-r}$ are fidelity parameters since the

simulator should not be expected to reach the true solution when some of these parameters are not at their highest fidelities, i.e., $\mathbf{t}_{-r} \neq \mathbf{0}$. The two limiting constraints thus describe how fidelity parameters determine the discrepancy behavior of the simulator: only when *all* fidelity parameters approach zero should the simulator converge to the true solution.

To satisfy these two properties, we place a Gaussian process prior on $\delta(\mathbf{x}, \mathbf{t})$ with product covariance form:

$$\text{Cov}[\delta(\mathbf{x_1}, \mathbf{t}_1), \delta(\mathbf{x_2}, \mathbf{t}_2)] = \sigma_2^2 K_{\mathbf{x}}^{\delta}(\mathbf{x}_1, \mathbf{x}_2) K_{\mathbf{t}}(\mathbf{t}_1, \mathbf{t}_2), \tag{13}$$

i.e., the effect of input variables and fidelity parameters are assumed to be separable for $\delta$. For the first kernel $K_{\mathbf{x}}^{\delta}(\cdot, \cdot)$, one can employ a standard stationary kernel; we make use of the squared-exponential kernel:

$$K_{\mathbf{x}}^{\delta}(\mathbf{x}_1, \mathbf{x}_2) = \exp \left\{ -\sum_{s=1}^{p} \alpha_s (x_{1,s} - x_{2,s})^2 \right\} \tag{14}$$

in our later implementation. For the second kernel $K_{\mathbf{t}}(\mathbf{t}_1, \mathbf{t}_2)$, one requires a careful *non-stationary* specification to satisfy the aforementioned two properties. Below, we present two choices for this kernel, which cater to two common scenarios for multi-stage multi-fidelity simulations:

- **Scenario 1**: The experiment consists of multiple fidelity parameters for simulating a *single mechanism* or phenomenon of interest. These different ways for varying accuracy, e.g., via spatial meshing and temporal discretization, represent different "stages" which control simulation precision. An example of this is the FEA of a cantilever beam deflection under stress, where three mesh fidelity parameters are used for each dimension of the three-dimensional finite element analysis. We will investigate this application further in Section 6.1.

- **Scenario 2**: The experiment comprises multiple stages performed *sequentially* over time, where a separate phenomenon is simulated at each stage. *Multiple mechanisms* are involved overall. Each stage may thus have various fidelity parameters controlling simulation precision. This is the case for our motivating nuclear physics problem

11

(Figure 1), where multiple mesh size parameters control simulation precision in each of the three consecutive stages for heavy-ion collisions; further details on this in Section 6.2.

In our experience, Scenario 2 is more often encountered than Scenario 1 in applications, but the kernel for Scenario 1 is perhaps more intuitive and is thus presented first. We note that these kernel choices are but recommendations – the modeler should carefully consider prior domain knowledge to carefully select a kernel that captures such knowledge. With the kernel $K_{\mathbf{t}}$ specified (along with kernels $K_{\mathbf{x}}^{\delta}$ and $K_{\mathbf{x}}^{\phi}$), one can show that the response surface $\eta(\mathbf{x}, \mathbf{t})$ follows a Gaussian process model, with covariance function:

$$K_{\eta}\{(\mathbf{x}_1, \mathbf{t}_1), (\mathbf{x}_2, \mathbf{t}_2)\} := \text{Cov}[\eta(\mathbf{x_1}, \mathbf{t_1}), \eta(\mathbf{x_2}, \mathbf{t_2})] = \sigma_1^2 K_{\mathbf{x}}^{\phi}(\mathbf{x}_1, \mathbf{x}_2) + \sigma_2^2 K_{\mathbf{x}}^{\delta}(\mathbf{x}_1, \mathbf{x}_2) K_{\mathbf{t}}(\mathbf{t}_1, \mathbf{t}_2). \quad (15)$$

The predictive equations for the M²GP then follow immediately from the standard GP equations (2) and (3) with kernel $K_{\eta}$ given above, with the desired prediction point $(\mathbf{x}^*, \mathbf{0})$ as the goal is predict the (limiting) highest-fidelity setting.

## 3.1   Kernel Option 1

Consider the first kernel choice for $K_{\mathbf{t}}$ (Kernel 1), which we recommend for Scenario 1 above. This takes the non-stationary form:

$$K_{\mathbf{t}}(\mathbf{t}_1, \mathbf{t}_2) = \exp\left\{-\sum_{r=1}^{k} \theta_r (t_{1,r} - t_{2,r})^2\right\} - \exp\left\{-\sum_{r=1}^{k} \theta_r t_{1,r}^2\right\} - \exp\left\{-\sum_{r=1}^{k} \theta_r t_{2,r}^2\right\} + 1. \quad (16)$$

Here, $\theta_r$ denotes the weight parameter for the $r^{\text{th}}$ fidelity parameter. A larger $\theta_r$ indicates greater sensitivity of discrepancy $\delta$ to the $r^{\text{th}}$ fidelity parameter, and vice versa. One can check that, with this kernel (16), the two desired properties (11) and (12) for $\delta$ are satisfied, meaning such a kernel indeed captures the aforementioned prior information on numerical convergence behavior.

Kernel 1 is inspired by the non-stationary covariance function in Gul et al. (2018), which was proposed for a different task of uncertainty propagation for system outputs.

It has several appealing features for multi-stage multi-fidelity emulation. First, in many applications, one may have prior knowledge of the *continuity* of the underlying numerical solutions (e.g., from FEA theory). With Kernel 1, the corresponding prior process on discrepancy $\delta$ can be shown to yield continuous sample paths, thus capturing such prior knowledge from a Bayesian perspective. Second, the form of this kernel provides a flexible framework for modeling *interactions* between fidelity parameters across different stages. Finally, the weight parameters $\theta_r$ of Kernel 1 provide flexibility in capturing the varying sensitivities of the simulator for different fidelity parameters, where both the differences between fidelity settings and their respective magnitudes are included. We later present both a maximum likelihood and Bayesian approach for inferring these parameters from data.

In our experience (see Section 6.1), Kernel 1 appears to work best for emulating a *single* mechanism with multiple fidelity parameters, e.g., the FEA for beam deflection with different fidelities for each of its three dimensions. One reason is that such systems often have significant interaction effects between fidelity parameters, e.g., between mesh sizes in each dimension, which appears to be well-captured by the non-stationary terms in (16).

## 3.2 Kernel Option 2

Consider next the second choice for $K_{\mathbf{t}}$ (Kernel 2), which we recommend for the multi-stage *sequential* simulations in Scenario 2. This kernel takes the non-stationary form:

$$K_{\mathbf{t}}(\mathbf{t}_1, \mathbf{t}_2) = \left[ \sum_{r=1}^{k} \theta_r \min(t_{1,r}, t_{2,r})^{l_r} \right]^{l}. \tag{17}$$

Here, $\theta_r$ is a weight parameter for the $r^{\text{th}}$ fidelity parameter, and $l_r$ and $l$ are kernel hyperparameters which we discuss later. Similar to Kernel 1, a greater $\theta_r$ allows for greater sensitivity of the discrepancy $\delta$ to the $r^{\text{th}}$ fidelity parameter. We can again show that with this kernel (17), the two properties (11) and (12) for bias $\delta$ are satisfied, meaning such a kernel also captures the desired prior information on numerical convergence. With Kernel 2, the resulting prior process on discrepancy $\delta$ can be viewed as a multivariate extension of

a standard Brownian motion model (Durrett 2019), and extends the non-stationary model (7) in Tuo et al. (2014), which tackled only the case of one fidelity parameter.

Kernel 2 has several appealing features for multi-stage multi-fidelity emulation, particularly when the multiple stages are performed *sequentially* over time (see Scenario 2 at the start of the section). One can show that the parametrization of this kernel is directly inspired by (and thus can capture prior information on) standard numerical convergence results for multi-stage simulators. To see why, consider first the simple setting of a *single* fidelity parameter $t$, and let $v_0$ and $v_t$ be the exact and simulated solutions respectively at fidelity $t$. In the case of finite element analysis (where $t$ is the mesh grid size), it is well-known (Brenner & Scott 2008) that the numerical error of the simulator can be upper bounded as:

$$\|\nu_0 - \nu_t\| \leq C t^\xi, \tag{18}$$

where $\|\cdot\|$ is an appropriate norm on the solution space, $\xi$ is a rate parameter, and $C$ is a constant. In words, the numerical error resulting from mesh discretization decays polynomially as mesh size $t$ decreases. Similar polynomial decay rates have also been shown for a broad range of fidelity parameters in numerical solvers, e.g., for elliptical PDEs (Hundsdorfer et al. 2003) and large-eddy simulations in fluid mechanics (Templeton et al. 2015).

Consider now the *multi-stage* simulations in Scenario 2, where a separate phenomenon is simulated sequentially at each stage. Suppose, at stage $r$, its precision is controlled by a fidelity parameter $t_r$. For this parameter $t_r$, further suppose the simulation error at this stage can be bounded by (18) with rate parameter $\xi_r$. One example of this is multi-stage finite element simulators when each stage involves a distinct finite element model (FEM) whose precision depends on a mesh size parameter $t_r$. Similar to before, let $\nu_\mathbf{0}$ and $\nu_{t_1,\cdots,t_k}$ denote the exact solution and the simulated solution at fidelity parameters $t_1, \cdots, t_k$. Applying the triangle inequality iteratively, the error between $\nu_{t_1,\cdots,t_k}$ and $\nu_\mathbf{0}$

can then be bounded as:

$$
\begin{aligned}
||v_{\mathbf{0}} - v_{t_1,\cdots,t_k}|| &\le ||v_{\mathbf{0}} - v_{t_1,0,\cdots,0}|| + ||v_{t_1,0,\cdots,0} - v_{t_1,t_2,0,\cdots,0}|| + \cdots \\
&\quad + ||v_{t_1,\cdots,t_{k-1},0} - v_{t_1,\cdots,t_{k-1},t_k}|| \\
&\le \sum_{r=1}^{k} C_r t_r^{\xi_r},
\end{aligned}
\tag{19}
$$

where $C_1, \cdots, C_k$ are again constants.

We now show that Kernel 2 indeed captures the error bound (19) as *prior information* within its kernel specification. To see why, consider the prior standard deviation of the discrepancy term $\delta(\mathbf{x}, \mathbf{t})$. From a Bayesian modeling perspective, this should capture the modeler's prior belief on the expected numerical error of the simulator. With $K_{\mathbf{t}}$ set as Kernel 2, one can show that this prior standard deviation takes the form:

$$
\sqrt{\operatorname{Var}\{\delta(\mathbf{x}, \mathbf{t})\}} = \sigma_2 \left[ \sum_{r=1}^{k} \theta_r t_r^{l_r} \right]^{l/2}.
\tag{20}
$$

Comparing (20) with (19), we see that they are precisely the same with the kernel hyperparameters set as $l = 2$ and $l_r = \xi_r$ for $r = 1, \cdots, k$. This suggests that, with $K_{\mathbf{t}}$ chosen as Kernel 2, the resulting prior model on discrepancy $\delta(\mathbf{x}, \mathbf{t})$ indeed captures (on expectation) the numerical error convergence of the multi-stage simulator.

The above connection also helps guide how the specification of hyperparameters for Kernel 2. If the rate parameters $\xi_1, \cdots, \xi_k$ can be identified via a careful analysis of the error bound (18) at each stage, one can use simply set the hyperparameters as $l_r = \xi_r$ for $r = 1, \cdots, k$. However, for more complex multi-stage simulators, one may not be able to identify the precise error convergence rates at each stage. In such cases, the kernel hyperparameters can be estimated via maximum likelihood or a fully Bayesian approach (see Section 4.1) or set at a fixed value (e.g., $l_r = l = 2$). Whether such hyperparameters are set a priori or inferred from data, the infusion of such prior information can yield noticeably improved predictive performance for multi-fidelity emulation, as we show later in Section 6.2.

15

# 4    Implementation

In this section, we discuss important implementation details for the M²GP model. We present two parameter inference approaches, the first via maximum likelihood and the second via a fully Bayesian formulation for incorporating external knowledge and richer uncertainty quantification. We then outline plausible experimental design strategies for the M²GP.

## 4.1    Parameter Inference

### 4.1.1    Maximum Likelihood

We first present a maximum likelihood approach for estimating the M²GP model parameters. Let $\boldsymbol{\Theta}_{\mathrm{MLE}} = (\boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\alpha}, \boldsymbol{\theta}, \sigma_1^2, \sigma_2^2)$ be the set of parameters to infer, where $\boldsymbol{\gamma}$ is the vector of weight parameters for $K_{\mathbf{x}}^{\phi}$, $\boldsymbol{\alpha}$ is the vector of weight parameters for $K_{\mathbf{x}}^{\delta}$, and $\boldsymbol{\theta}$ is the vector of weight parameters for the M²GP kernel $K_{\mathbf{t}}$ (either Kernel 1 or Kernel 2). Since $\eta(\mathbf{x}, \mathbf{t})$ can be expressed as a GP with kernel given in (15), one can easily obtain an analytic expression for the likelihood function to optimize. More specifically, let the simulated multi-fidelity training data be $\mathbf{y} = (\eta(\mathbf{x}_i, \mathbf{t}_i))_{i=1}^n$, and let the matrix of basis functions for the GP mean be $\mathbf{F} = (\mathbf{f}(\mathbf{x}_1, \mathbf{t}_1)^T; \mathbf{f}(\mathbf{x}_2, \mathbf{t}_2)^T; \cdots ; \mathbf{f}(\mathbf{x}_n, \mathbf{t}_n)^T)$, with corresponding coefficients $\boldsymbol{\beta}$. We thus aim to maximize the log-likelihood of the M²GP model, given by:

$$\max_{\boldsymbol{\Theta}_{\mathrm{MLE}}} \left\{ -\frac{1}{2} \log \det \boldsymbol{\Sigma} - \frac{1}{2} (\mathbf{y} - \mathbf{F}\boldsymbol{\beta})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{F}\boldsymbol{\beta}) \right\}, \tag{21}$$

where $\det \boldsymbol{\Sigma}$ is the determinant of the covariance matrix $\boldsymbol{\Sigma}$.

While the optimization problem (21) is quite high-dimensional, standard non-linear optimization algorithms, such as the L-BFGS-B method (Nocedal & Wright 1999), appear to work reasonably well. One can further speed up this optimization procedure via an informed initialization of the parameters $\boldsymbol{\Theta}_{\mathrm{MLE}}$. In particular, we have found that the correlation parameters $\boldsymbol{\alpha}$ can be well-initialized by first fitting a standard Gaussian process model on $\mathbf{y}$ with kernel $K_{\mathbf{x}}^{\delta}$. With these initial estimates, we then perform the L-BFGS-B

non-linear optimization algorithm, as implemented in the R package `stats` (Byrd et al. 1995).

### 4.1.2 Fully Bayesian Inference

In situations where a richer quantification of uncertainty is desired, a fully Bayesian approach to parameter inference may be appropriate. Below, we present one such approach for the $\mathrm{M}^2\mathrm{GP}$ which leverages a Metropolis-within-Gibbs algorithm (Gelman et al. 1995) for posterior sampling. For easier derivation of the full conditional distributions, we consider a reparametrization of the covariance kernel (15) for $\eta(\mathbf{x}, \mathbf{t})$ as:

$$K_\eta\{(\mathbf{x}_1, \mathbf{t}_1), (\mathbf{x}_2, \mathbf{t}_2)\} = \sigma^2 \left\{ K_{\mathbf{x}}^\phi(\mathbf{x}_1, \mathbf{x}_2) + \lambda K_{\mathbf{x}}^\delta(\mathbf{x}_1, \mathbf{x}_2) K_{\mathbf{t}}(\mathbf{t}_1, \mathbf{t}_2) \right\}, \tag{22}$$

where $\sigma^2 \coloneqq \sigma_1^2$ and $\lambda \coloneqq \sigma_2^2/\sigma_1^2$. Here, the new parameter $\lambda$ captures the degree of non-stationarity in the kernel from the influence of the fidelity parameters $\mathbf{t}$. When $\lambda = 0$, the covariance kernel becomes a stationary kernel that depends on only input parameters $\mathbf{x}$.

With this reparametrization, the parameter set to infer is given by $\boldsymbol{\Theta}_\mathrm{B} = (\boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\alpha}, \boldsymbol{\theta}, \sigma^2, \lambda)$. It is straightforward to show that:

$$\mathbf{y}|\boldsymbol{\Theta}_\mathrm{B} \sim \mathcal{N}(\mathbf{F}\boldsymbol{\beta}, \boldsymbol{\Sigma}), \tag{23}$$

using the same notation as in (21). Table 1 summarizes the priors assigned on parameters $\boldsymbol{\Theta}_\mathrm{B}$. Here, the prior hyperparameters can either be set via prior information, or set in a weakly-informative fashion with $a_\lambda = b_\lambda = 1$ and $a = b = 0.001$ for the remaining hyperparameters.

With the priors specified, we now proceed to the posterior sampling algorithm. Of the model parameters in $\boldsymbol{\Theta}_\mathrm{B}$, we can derive full conditional distributions for two parameters, $\boldsymbol{\beta}$ and $1/\sigma^2$:

$$\boldsymbol{\beta}|\mathbf{y}, \boldsymbol{\gamma}, \boldsymbol{\alpha}, \boldsymbol{\theta}, \sigma^2, \lambda \sim \mathcal{N}((\mathbf{F}^T\boldsymbol{\Sigma}^{-1}\mathbf{F})^{-1}\mathbf{F}^T\boldsymbol{\Sigma}^{-1}\mathbf{y}, \sigma^2(\mathbf{F}^T\boldsymbol{\Sigma}^{-1}\mathbf{F})^{-1}), \tag{24}$$

$$1/\sigma^2|\mathbf{y}, \boldsymbol{\gamma}, \boldsymbol{\alpha}, \boldsymbol{\theta}, \lambda, \boldsymbol{\beta} \sim \mathrm{Gamma}\left(a_\sigma + \frac{n}{2}, (1+\lambda)b_\sigma + \frac{1}{2}(\mathbf{y} - \mathbf{F}\boldsymbol{\beta})^T\boldsymbol{\Sigma}^{-1}(\mathbf{y} - \mathbf{F}\boldsymbol{\beta})\right). \tag{25}$$

| Model | Prior Specification |
|-------|---------------------|
| M²GP:    $\eta(\mathbf{x}, \mathbf{t}) \sim \mathcal{GP}\{\boldsymbol{F}\boldsymbol{\beta}, K_\eta(\cdot, \cdot)\}$ | $[\beta_1, \beta_2, \cdots, \beta_m] \overset{i.i.d.}{\sim} 1$ |
| Priors:    $[\boldsymbol{\Theta}_{\mathrm{B}}] = [\boldsymbol{\beta}][\lambda][\sigma^2|\lambda][\boldsymbol{\gamma}][\boldsymbol{\alpha}][\boldsymbol{\theta}]$ | |
| *Non-stationary parameter* | $\lambda \sim \mathrm{Beta}(a_\lambda, b_\lambda)$ |
| *Kernel precision* | $1/\sigma^2|\lambda \sim \mathrm{Gamma}(a_\sigma, (1+\lambda)b_\sigma)$ |
| *Weight parameters* | $\gamma_1, \gamma_2, \cdots, \gamma_p \overset{i.i.d.}{\sim} \mathrm{Gamma}(a_\gamma, b_\gamma)$ |
| *Weight parameters* | $\alpha_1, \alpha_2, \cdots, \alpha_p \overset{i.i.d.}{\sim} \mathrm{Gamma}(a_\alpha, b_\alpha)$ |
| *M²GP weight parameters* | $\begin{cases} \theta_1, \theta_2, \cdots, \theta_k \overset{i.i.d.}{\sim} \mathrm{Gamma}(a_\theta, b_\theta) \text{ for Kernel } 1 \\ \theta_1, \theta_2, \cdots, \theta_k \overset{i.i.d.}{\sim} \mathrm{Beta}(a_\theta, b_\theta) \text{ for Kernel } 2 \end{cases}$ |

Table 1: Hierarchical model specification for the fully Bayesian M²GP.

For the remaining parameters in $\boldsymbol{\Theta}_{\mathrm{B}}$, we make use of Metropolis-Hastings (Metropolis et al. 1953) steps for sampling the full conditional distributions, as implemented in the R package `MHadaptive` (Chivers 2012). We then iterate these full conditional sampling steps within a Gibbs sampler for posterior exploration of $[\boldsymbol{\Theta}_{\mathrm{B}}|\mathbf{y}]$. Algorithm 1 presents the detailed steps for this Metropolis-within-Gibbs sampler for the M²GP.

Finally, with the posterior samples $\{\boldsymbol{\Theta}_{\mathrm{B}}^{[m]}\}_{m=1}^M$ obtained from Algorithm 1, we can easily estimate the posterior predictive mean at a new test point $\mathbf{x}^*$ by marginalizing over $\boldsymbol{\Theta}_{\mathrm{B}}$:

$$\mathbb{E}[\eta(\mathbf{x}^*, \mathbf{0})|\mathbf{y}] \approx \frac{1}{M} \sum_{m=1}^M \hat{\eta}(\mathbf{x}^*, \mathbf{0}|\boldsymbol{\Theta}_{\mathrm{B}}^{[m]}),$$

where $\hat{\eta}(\mathbf{x}^*, \mathbf{0}|\boldsymbol{\Theta}_{\mathrm{B}}^{[m]})$ is the closed-form GP predictive mean in (2) with fixed hyperparameters $\boldsymbol{\Theta}_{\mathrm{B}}^{[m]}$. This serves as the emulator for the fully-Bayesian M²GP. One can also quantify its uncertainty via posterior predictive samples on $\eta(\mathbf{x}^*, \mathbf{0})|\mathbf{y}$. These can be obtained by sampling a batch of samples from the predictive distribution $[\eta(\mathbf{x}^*, \mathbf{0})|\mathbf{y}, \boldsymbol{\Theta}_{\mathrm{B}}^{[m]}]$ in (2) given parameters $\boldsymbol{\Theta}_{\mathrm{B}}^{[m]}$, then repeating this procedure on all posterior samples $\{\boldsymbol{\Theta}_{\mathrm{B}}^{[m]}\}_{m=1}^M$.

**Algorithm 1** Metropolis-within-Gibbs sampler for the $\mathrm{M}^2\mathrm{GP}$

**Input:** Training data $\{\mathbf{x}_i, \mathbf{t}_i\}_{i=1}^n$, $\mathbf{y} = (\eta(\mathbf{x}_i, \mathbf{t}_i))_{i=1}^n$; testing input $\mathbf{x}^*$; prior hyperparameters $a_\lambda$, $b_\lambda$, $a_\sigma$, $b_\sigma$, $a_\gamma$, $b_\gamma$, $a_\alpha$, $b_\alpha$, $a_\theta$, $b_\theta$; number of MCMC samples $M$.

**Output:** Samples from the posterior distribution $[\mathbf{\Theta}_\mathrm{B}|\mathbf{y}]$.

1: Initialize the parameters $\mathbf{\Theta}_\mathrm{B}^{[0]}$ from the prior.

2: **for** $m = 1, \cdots, M$ **do**

3:     Sample $\boldsymbol{\beta}^{[m]}$ from the full conditional distribution (24).

4:     Sample $1/\sigma^{2[m]}$ from the full conditional distribution (25).

5:     For the remaining parameters $\{\boldsymbol{\gamma}, \boldsymbol{\alpha}, \boldsymbol{\theta}, \boldsymbol{\lambda}\}$, perform one step of Metropolis-within-Gibbs sampling using parameters $\{\boldsymbol{\beta}^{[m]}, 1/\sigma^{2[m]}\}$.

6: **end for**

**Return:** Posterior samples $\{\mathbf{\Theta}_\mathrm{B}^{[m]}\}_{m=1}^M$ from the posterior distribution $[\mathbf{\Theta}_\mathrm{B}|\mathbf{y}]$.

## 4.2 Experimental Design

Of course, given a fixed and limited computational budget, an experimenter would want to maximize the predictive power of the fitted multi-fidelity emulator model. For Gaussian process models, space-filling designs (Joseph 2016) – which aim to uniformly fill up the design space – are commonly used, and have desirable information-theoretic and predictive properties (Johnson et al. 1990). Different notions of space-filling designs have been explored in the literature, including maximin designs (Johnson et al. 1990, Morris & Mitchell 1995), minimax designs (Johnson et al. 1990, Mak & Joseph 2018) and maximum projection (MaxPro) designs (Joseph et al. 2015).
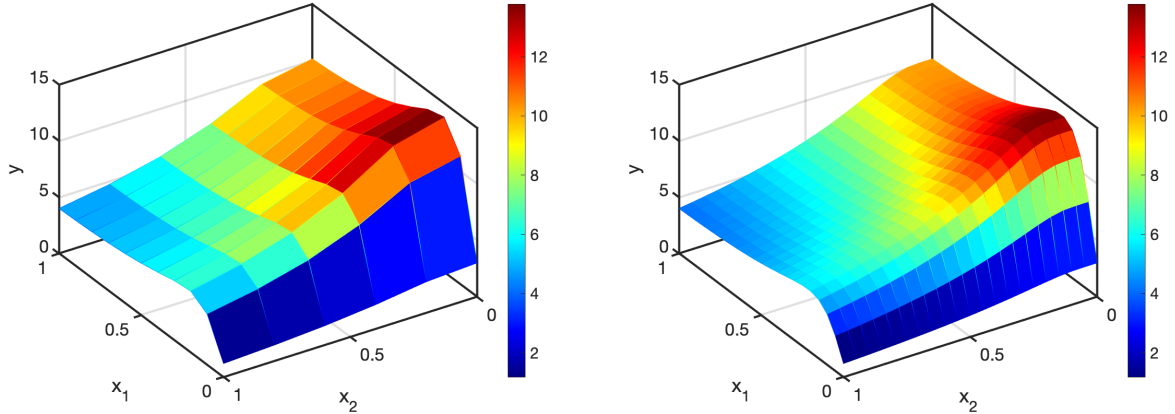
For the $\mathrm{M}^2\mathrm{GP}$ model, there are several ways in which one can adapt existing space-filling designing methods. One approach is to (i) adopt a space-filling design over the combined design space of both input parameters $\mathbf{x}$ and fidelity parameters $\mathbf{t}$. Such a design ensures that training points are not only well-spaced out over the input space for prediction at untested settings but also well-spaced out over the fidelity space to better learn the effects of individual fidelity parameters. Another approach might be (ii) a crossed array design

(Wu & Hamada 2009) between input and fidelity space, which are popular designs for robust parameter design. In such a design, one first generates two space-filling designs, one over the input space and the other over the fidelity space, then takes for the final design all combinations of input and output points. Both designs appear to yield good performance: the designs in (i) are used for our numerical experiments and cantilever beam deflection application, and the designs in (ii) are used for the emulation of the QGP evolution. The problem of optimal experimental design for the proposed non-stationary $\text{M}^2\text{GP}$ model is quite intriguing, and we aim to pursue this in future work.

# 5    Numerical Experiments

We now explore the performance of the proposed $\text{M}^2\text{GP}$ multi-fidelity model in a suite of simulation experiments with multiple fidelity parameters. We compare the $\text{M}^2\text{GP}$ with several existing emulator models. The first model is a standard GP emulator with a squared-exponential correlation function on both input parameters $\mathbf{x}$ and fidelity parameters $\mathbf{t}$; one then uses the fitted model to predict at $\mathbf{t} = \mathbf{0}$. We call this model simply the "standard GP" emulator. The second model is the TWY model (Tuo et al. 2014), which uses a single fidelity parameter. Since there are multiple fidelity parameters in the target multi-stage problem, we will first compute the arithmetic or geometric mean of the fidelity parameters $t_1, \cdots, t_k$, then apply the TWY model with this single aggregate fidelity parameter. We call the resulting models the TWY (ARITH) and TWY (GEOM) emulators, respectively.

In the following, we investigate the performance of these models on multi-fidelity extensions of two test functions, the 2D Currin function (Currin et al. 1991) and the 4D Park function (Cox et al. 2001). For the $\text{M}^2\text{GP}$, we follow Section 4 and set the power parameters in Kernel 2 as $l_r = l = 2$. Kernel hyperparameters for our model are estimated via maximum likelihood in Sections 5.1 and 5.2, and its fully Bayesian counterpart is explored in Section 5.3.

(a) $\eta(\mathbf{x}, \mathbf{t})$ with $t_1 = 0.1$ and $t_2 = 0.2$.  (b) $\eta(\mathbf{x}, \mathbf{t})$ with $t_1 = 0.05$ and $t_2 = 0.05$.

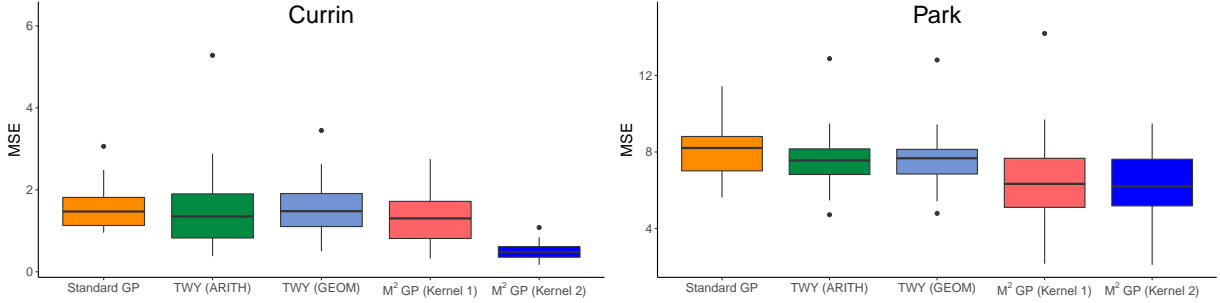Figure 2: Visualization of the multi-fidelity Currin function at two different fidelity settings.

## 5.1  Multi-Fidelity Currin Function

Our first test function builds off of the 2D Currin test function in Currin et al. (1991):

$$\phi(\mathbf{x}) = \left[1 - \exp\left(-\frac{1}{2x_2}\right)\right] \frac{2300x_1^3 + 1900x_1^2 + 2092x_1 + 60}{100x_1^3 + 500x_1^2 + 4x_1 + 20}, \tag{26}$$

where $\mathbf{x} = [x_1, x_2] \in [0,1]^2$. We then build a lower-fidelity representation of this function, denoted as $\eta(\mathbf{x}, \mathbf{t})$, with two fidelity parameters $\mathbf{t} = [t_1, t_2]$ via piecewise grid interpolation. More specifically, this approximation is carried out in two steps. First, we generate a rectangular grid in the input space, where the dimension of each mesh cell is $t_1 \times t_2$. Next, we evaluate the underlying function (26) at the mesh grid points, and perform piecewise grid interpolation to construct a lower-fidelity version of (26). This procedure is effectively the same as finite element meshing, which splits the input domain into many smaller elements. Fig. 2 visualizes this test function $\eta(\mathbf{x}, \mathbf{t})$ with $(t_1, t_2) = (0.1, 0.2)$ and $(0.05, 0.05)$. It is clear that, as $t_1$ and $t_2$ become smaller, $\eta(\mathbf{x}, \mathbf{t})$ becomes closer to the underlying Currin function (26), which is as desired.

We then compare the M²GP emulator with the aforementioned baseline models. For each experiment, we first generate $n = 50$ design points over *both* input and fidelity parameters via the MaxPro design (Joseph et al. 2015). Here, we set the range for each

| Model | Average MSE (Currin) | Average MSE (Park) |
|-------|----------------------|--------------------|
| Standard GP | 1.519 | 8.132 |
| TWY (ARITH) | 1.525 | 7.580 |
| TWY (GEOM) | 1.560 | 7.566 |
| M²GP (Kernel 1) | **1.318** | **6.537** |
| M²GP (Kernel 2) | **0.509** | **6.289** |

Figure 3: (Top left) Boxplots of MSEs for the multi-fidelity Currin function experiment. (Top right) Boxplots of MSEs for the multi-fidelity Park function experiment. (Bottom) Average MSEs for the multi-fidelity Currin and Park experiments.

fidelity parameter to be between 0.1 and 0.4, to mimic the reality that simulations are prohibitively expensive for small choices of fidelity parameters $t_i$. Using this design, we then collect training data from the multi-fidelity Currin function $\eta(\mathbf{x}, \mathbf{t})$. For validation, we randomly select $N = 1,000$ points over the input parameter space as the testing set and compare how well these models predict the desired Currin function $\phi(\mathbf{x})$ via mean squared error (MSE). This procedure is then replicated 30 times.

Figure 3 (top left) shows the boxplots of the testing MSE for the proposed M²GP model with Kernel 1 and 2, as well as the standard GP emulator and the two variants of the TWY model, with Figure 3 (bottom) summarizing their average MSEs. There are several interesting observations to note. First and foremost, the M²GP emulator (with either Kernel 1 or Kernel 2) noticeably outperforms the existing emulator models. This suggests

that, by embedding prior information on the known convergence of $\eta(\mathbf{x}, \mathbf{t})$ within the non-stationary kernel specification, the proposed M²GP model can indeed provide improved emulation performance over models which do not explicitly integrate such information. The advantages of integrating the *multi-stage* multi-fidelity framework can be seen by comparing the M²GP models with the TWY models. By directly modeling the effects of multiple fidelity parameters, the proposed models provide a more faithful representation of the multi-stage simulator and thus improved predictions over the TWY models, which aggregate fidelity into a single parameter. Finally, the M²GP with Kernel 2 provides noticeably better performance than Kernel 1; this may be because there is little interaction between the two fidelity parameters under the piecewise grid interpolation of $\eta(\mathbf{x}, \mathbf{t})$.

## 5.2 Multi-Fidelity Park Function

Our second test function builds off of the 4D Park test function in Cox et al. (2001):

$$\phi(\mathbf{x}) = \frac{x_1}{2} \left[ \sqrt{1 + (x_2 + x_3^2)\frac{x_4}{x_1^2}} - 1 \right] + (x_1 + 3x_4) \exp\left(1 + \sin(x_3)\right), \tag{27}$$

We again build a lower-fidelity representation of $\phi(\mathbf{x})$, denoted as $\eta(\mathbf{x}, \mathbf{t})$, using four fidelity parameters $\mathbf{t} = (t_1, \cdots, t_4)$ via piecewise grid interpolation. Similar to before, we use MaxPro designs (with $n = 50$ design points) over both input and fidelity parameters, with a range of $[0.2, 0.5]$ for each fidelity parameter. The same emulator models are compared as before, and the experiment is replicated 20 times over $N = 1,000$ random testing points.

Figure 3 (top right) shows the boxplots of the testing MSE for the five models, with Figure 3 (bottom) reporting its average MSEs. We see that the proposed M²GP model again outperforms its competitors by a noticeable margin, which affirms the value of embedding prior information on the multi-stage convergence of $\eta(\mathbf{x}, \mathbf{t})$ within the M²GP non-stationary kernel specification. Similarly, Kernel 2 provides slightly improved performance over Kernel 1, which may be due to little interaction between fidelity parameters for piecewise grid interpolation.
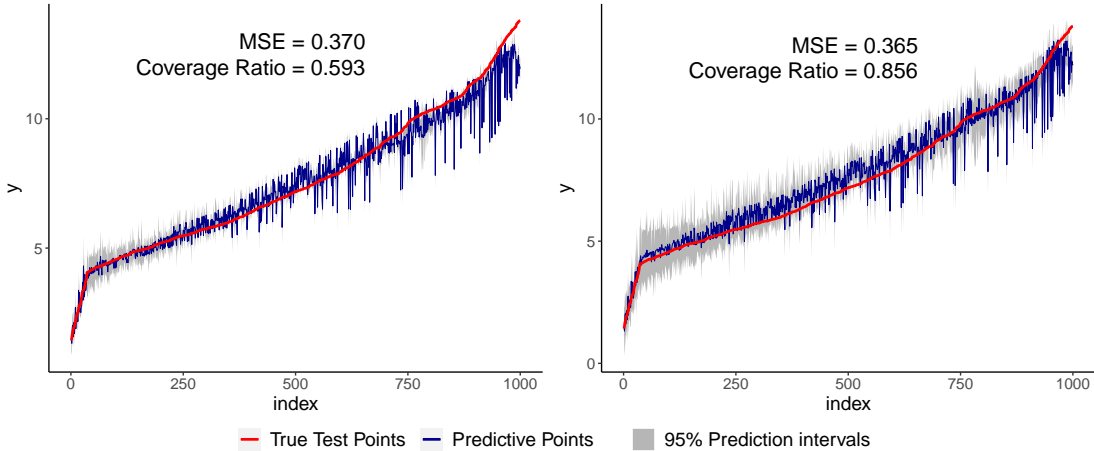
Figure 4: Visualizing the predictive performance for the M²GP model with a plug-in MLE (left) and fully Bayesian (right) implementation for the multi-fidelity Currin simulation. Here, the true function values are plotted in red, emulator predictions in blue, and 95% predictive intervals in gray. Also reported are the testing MSEs and coverage ratios.

## 5.3    Fully Bayesian Implementation

In previous experiments, the M²GP model parameters were estimated via maximum likelihood, which may lead to underestimation of predictive uncertainty; when this is the case, a fully Bayesian implementation (see Section 4.1.2) may be preferable. Here, we compare the performance of both implementations of the M²GP for the earlier multi-fidelity Currin function experiment. Both approaches are compared in terms of their predictive accuracy (i.e., its testing MSE) and coverage ratio (how well its 95% predictive intervals cover the test points). For the MLE approach, its 95% predictive interval is obtained from the closed-form distribution (2) with plug-in parameter estimates. For the fully-Bayesian approach, its 95% predictive interval is computed from posterior samples on the predictive distribution $[\eta(\mathbf{x}^*, \mathbf{0})|\mathbf{y}]$, with the MCMC chain run for 50,000 iterations and thinned to 100 samples. Since the fully Bayesian model is more expensive to fit, we use only the first set of training/testing data from the earlier experiment. This comparison is made using Kernel 2, which has the lowest MSE from Section 5.1 for the Currin function.

Figure 4 summarizes the MSEs and coverage ratios over the test set, for the MLE and fully Bayesian implementations of the M²GP. We see that, while the testing MSE of the MLE approach is quite small, its coverage ratio is around 60%, which is significantly lower than the nominal rate of 95%. This is not surprising: not only is the quantification of uncertainty for *extrapolation* problems challenging (since prediction at $t_1 = t_2 = 0$ is outside of the range for training data), it is further hampered by the lack of parameter uncertainty captured by the plug-in MLE approach. The fully Bayesian approach yields similar testing MSE to the MLE approach but provides noticeably closer coverage (around 86%) to the desired rate of 95%. While we still get a slight underestimation of uncertainty (which is again unsurprising since extrapolation using GPs is inherently difficult), the integration of parameter uncertainty via a hierarchical Bayesian framework can indeed provide for a richer quantification of uncertainty for emulation.

# 6 Applications

Finally, we explore the usefulness of the proposed model in two applications. The first application involves the multi-stage multi-fidelity emulation of a cantilever beam deflecting under stress. The second application is the earlier motivating problem of multi-stage multi-fidelity emulation of the quark-gluon plasma produced in heavy-ion collisions.

## 6.1 Cantilever Beam Deflection

The first application investigates the static stress analysis on a cantilever beam. Beam structures are commonly used in finite elements to model transverse loads and deformation under various circumstances, and the study of their deflection behavior is a canonical problem in finite element analysis and has been studied extensively (Ngo & Scordelis 1967, Heyliger & Reddy 1988, Chakraborty et al. 2003). Here we use it to evaluate the performance of our modeling framework. Figure 5 shows an illustration of the beam for our study, where one end surface of the beam is fixed, and an external pressure field is applied on
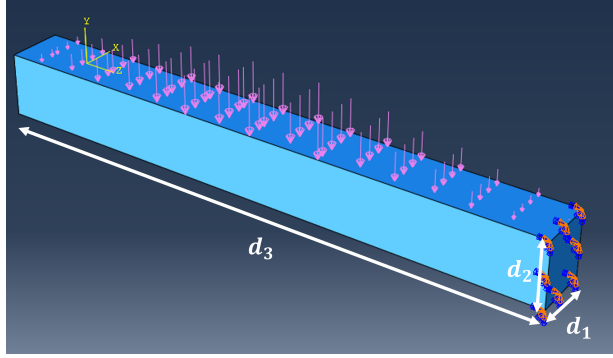
Figure 5: Beam cantilever simulation with fixed end surface as shown in ABAQUS.

its top surface. The deflection of this beam under stress is typically simulated using FEA simulations, which can be computationally expensive. For our experiments, these FEA simulations are carried out using the ABAQUS software (Smith 2009) with rectangular mesh cells.

The set-up is as follows. The beam dimensions are specified by its breadth $d_1$, its height $d_2$, and its length $d_3$. We further let $d_1 = d_2$ so the cross section of the beam is square-shaped (see Figure 5). We then set the Young's modulus of the beam (which parametrizes the stiffness of the beam) to be 200 MPa and the Poisson ratio (which measures the deformation of the beam under loading) to be 0.28, with material properties corresponding to steel. For the external pressure field, which is applied vertically downward on top of the beam, we employed the continuous half-sine pressure field given by:

$$p(l) = C_1 x_1 \sin(C_2 l / x_2). \tag{28}$$

Here, $l \in [0, d_3]$ denotes the location along the beam from the fixed end, $C_1 = 2000$ and $C_2 = \pi/200$ are constants, $x_1$ is a scale factor for the pressure, and $x_2$ parametrizes the length of the beam, i.e., $d_3 = 200x_2$. An additional input parameter $x_3$ controls the width and breath of the beam cross section $x_3 = 20d_1 = 20d_2$. There are thus a total of three input parameters $\mathbf{x} = [x_1, x_2, x_3] \in [0, 1]^3$ for this study.

For fidelity parameters, it is natural to consider a meshing procedure which partitions the beam into smaller 3D mesh rectangles. The size of these mesh rectangles can be

controlled by three fidelity parameters, which dictate the size of the mesh rectangles in each dimension. In other words, the three fidelity parameters $t_1, t_2, t_3 \in (0, 1)$ determine the *scale* of the finite elements. As a result of the pressure field, the cantilever beam will deflect downward, resulting in deflection at its tip. The response of interest is taken to be the amount of tip deflection. The goal is thus to train an emulator model which, using a carefully designed training set of simulation runs over different inputs $\mathbf{x}$ and fidelities $\mathbf{t}$, efficiently predicts the "exact" solution for tip deflection (i.e., at $\mathbf{t} = \mathbf{0}$) of a new beam with inputs $\mathbf{x}$.

The experiment is carried out as follows. To generate training data, we first run the simulator (ABAQUS) on a $n = 50$-point MaxPro design (Joseph et al. 2015) over the combined space of input parameters $\mathbf{x}$ and fidelity parameters $\mathbf{t}$, which required about 4.5 hours of computation. For fidelity parameters, we set it to be $\mathbf{t} \in [1/31, 1/3]^3$, which ensures we have an integer number of finite elements at the edge case in each dimension (for $\mathbf{t}$ values in between, we round up to the nearest integer). For validation, we further run the simulator on 30 new cantilever configurations (the testing set) where each takes about one hour, uniformly sampled over the input space, to test the performance of each model (in terms of mean absolute error) in predicting the tip deflections. While the "exact" response with $\mathbf{t} = \mathbf{0}$ cannot be obtained numerically, this can be well-approximated by running the simulator at very fine mesh sizes; in our case, we used $\mathbf{t} = [0.025, 0.025, 0.005]$ for testing points, which provided a sufficiently fine mesh according to a mesh validation study. One simulation run at this high-fidelity setting requires around an hour of computation, meaning there is a considerable opportunity for a multi-fidelity emulator to greatly speed up design exploration. The same emulators as before (the standard GP, the two TWY models, and the two M²GP models) are used for comparison here. In addition, we include a "high-fidelity GP" emulator model, which is trained on only data from the high-fidelity simulator with $\mathbf{t} = [0.025, 0.025, 0.005]$. For a fair comparison, this model is trained on high-fidelity points from a four-point MaxPro design, which requires comparable time to simulate as the earlier 50-point designs over the combined input-fidelity space.
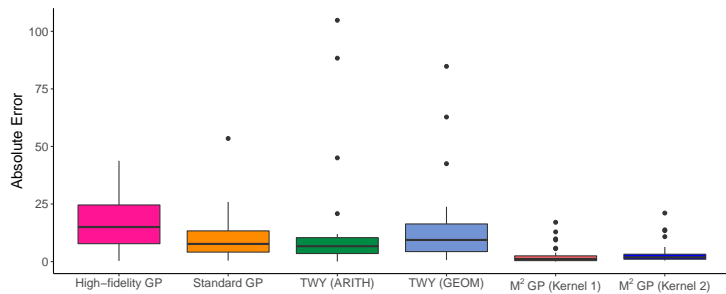
Figure 6: Boxplots of absolute errors for the beam deflection application.

| Model | MAE |
|---|---|
| High-fidelity GP | 17.52 |
| Standard GP | 10.32 |
| TWY (ARITH) | 13.52 |
| TWY (GEOM) | 14.75 |
| $M^2GP$ (Kernel 1) | **2.83** |
| $M^2GP$ (Kernel 2) | **3.65** |

Table 2: MAEs for the beam deflection application.

Table 2 summarizes the mean absolute errors (MAEs) for the emulators compared in this application. Again, we see that the proposed $M^2GP$ model (with either Kernel 1 or 2) yields noticeably improved predictive performance over existing methods, with MAEs roughly 68% lower than for the standard GP, and more than 80% lower than the TWY models. We also see that, given a similar computational budget for training data simulation, the high-fidelity emulator yields noticeably poorer performance compared to the multi-fidelity emulators, which is unsurprising. Figure 6 shows the boxplots of absolute prediction errors, which confirm the earlier observation. It is interesting to note that for the $M^2GP$, Kernel 1 provides slightly better predictive performance compared to Kernel 2. One reason is that, as a single-mechanism multi-fidelity problem, this cantilever beam deflection application can be classified under Scenario 1 (see Section 3), where the experiment simulates a *single* mechanism with multiple fidelity parameters. Kernel 1 appears to be better suited at capturing the more complex interactions between fidelity parameters, thus leading to slightly better performance over Kernel 2.

Here, we again see that by integrating prior information on the numerical convergence of the multi-stage multi-fidelity simulator within its non-stationary kernel specification, the proposed $M^2GP$ can yield noticeably improved emulation performance over existing methods. This is particularly apparent in the cantilever beam application, where the
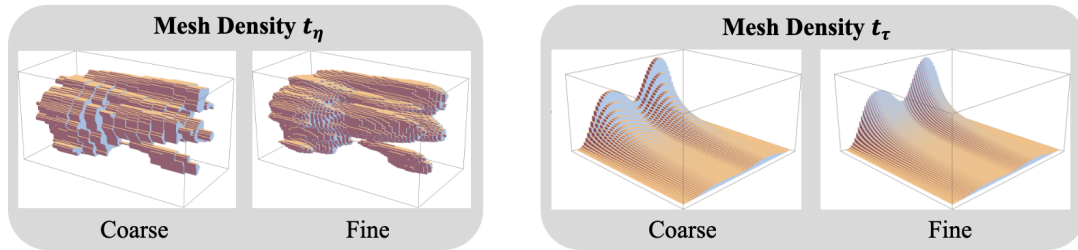
28

Figure 7: Visualizing the two mesh densities (fidelity parameters) in the quark-gluon plasma simulation.

three fidelity parameters bear physical importance. For bending beams, the accuracy of simulations is generally considered to be more sensitive to mesh density along the span of the beam (i.e., $d_3$) than along the other two dimensions (Cui & Wisnom 1992). By explicitly accounting for this multi-parameter fidelity structure, the M²GP can pick out the greater importance of fidelity parameter $t_3$ via inference of its weight parameters, thus allowing for improved predictive accuracy over existing models which ignore such structure.

## 6.2 Quark-Gluon Plasma Evolution

We now return to our motivating problem of emulating the quark-gluon plasma, an exotic state of nuclear matter which can be created in modern particle colliders, and which pervaded the universe during its first microseconds. The study of this plasma – in particular, properties of this unique phase of matter – is thus an important problem in high-energy nuclear physics. Modern investigation of QGP often requires computationally-demanding numerical simulations, with the plasma modeled via relativistic fluid dynamics. The use of cost-efficient emulators, when carefully constructed, can thus greatly speed up the discovery of fundamental properties on the QGP, as evidenced in recent works (Liyanage et al. 2022, Cao et al. 2021).

For this study, we adopt a simplified version of the QGP simulation framework in Everett et al. (2021), which can be split into three distinct stages: a pre-hydrodynamic stage, a hydrodynamic stage, and a post-hydrodynamic stage. Figure 1 visualizes this multi-stage

simulation framework. Each stage typically involves discretization of the simulated physical system onto a spatial or space-time mesh (see Figure 7). The sizes and dimensionalities of the meshes may vary among stages. Meshes must be large enough to contain the entire initial and final states of the systems, fine enough to capture relevant details (e.g., small-scale fluctuations in the pre-hydrodynamic initial state), and yet allow timely computation.

The considered simulator has two key fidelity parameters, $t_\eta$ and $t_\tau$, which control its precision. The first fidelity parameter arises in the pre-hydrodynamic stage. This stage models the initial distribution of energy resulting from the collision of two atomic nuclei. The energy distribution is defined on a 3D (spatial) mesh with coordinates $x$, $y$, and $\eta$. The bounds of the mesh are fixed in all three dimensions, but the mesh density in the $\eta$ direction will be varied to adjust fidelity; it is specified by the longitudinal mesh size variable $t_\eta$, which serves as our first fidelity parameter. The simulation costs of all three stages are inversely proportional to $t_\eta$. The second parameter arises when the initial hydrodynamic state is evolved with the relativistic hydrodynamic equations until a completion criterion is reached, in effect extending the mesh into a time dimension, denoted $\tau$. The temporal mesh size variable $t_\tau$ – our second fidelity parameter – can thus be varied to adjust fidelity, although we note that in contrast with the $\eta$ spatial direction, the number of timesteps is not known in advance because the full evolution time cannot be fixed – it is only determined once the completion criterion is satisfied, and so depends on the initial conditions in a complicated way. Except at very low fidelity, the simulation costs of the hydrodynamic and post-hydrodynamic stages are inversely proportional to $t_\tau$.

In this simplified QGP simulator, we consider a single response variable: the ratio of pions produced at two different points, $\eta = 0$ and $\eta = 1$. This ratio serves as a measure of how particle production is distributed along the collision axis of the atomic nuclei, and is chosen because it is strongly influenced by the model parameter $\alpha$, which we use as our single input parameter in this study. We denote $\alpha$ as $x_1$ and the ratio observable as $y_1$ below.

The experiment is carried out as follows. We compare the M$^2$GP with the standard

GP model and the TWY (ARITH) model; the TWY (GEOM) model was very numerically unstable and was excluded in this comparison. Since the multi-stage procedure involves multiple sequential stages, it falls under Scenario 2 (see Section 3), and thus we make use of Kernel 2 for M$^2$GP. To demonstrate the cost efficiency of multi-fidelity emulation, we again include the "high-fidelity GP" model, which makes use of *only* high-fidelity runs to train a standard GP emulator using the squared-exponential kernel. As before, the limiting highest-fidelity setting of $\mathbf{t} = \mathbf{0}$ cannot be numerically simulated. We thus set the fidelity parameters $\mathbf{t} = (1.0 \times 10^{-4}, 1/64)$ as the "high-fidelity" setting for prediction, which appears to provide a fine enough mesh according to a mesh validation study. With this, a single high-fidelity run is very time-consuming, requiring around 1,000 CPU hours.

For comprehensive cost analysis, we fit each emulator using different design sizes, then compare the predictive performance of these models given a computational budget. The training data are generated as follows. For the high-fidelity GP model, we generate $n = 2, 3, 4,$ or $5$ maximin (equally-spaced) high-fidelity design points over the input interval $x_1 \in [3, 5]$. For the remaining models, we generate $n = 15, 20,$ or $25$ design points over the joint space of input and fidelity parameters. Each design has an equal number of points on five maximin (equally-spaced) levels on $x_1$. For the two fidelity parameters $t_\tau$ and $t_\eta$, we first generate a 2D $n$-point maximin LHD (Morris & Mitchell 1995) and scale this over the domain $[1.0 \times 10^{-4}, 5.0 \times 10^{-2}] \times [1/64, 1/24]$. We then randomly assign to each level of $x_1$ a fidelity setting from this LHD. For validation, the test set is generated on 100 evenly-spaced points over the input space $x_1 \in [3, 5]$, run at the aforementioned high-fidelity setting for $\mathbf{t}$. This procedure is then replicated 20 times. The median computational cost for the training data ranges from $1.5 \times 10^3$ (15 points) to $3.3 \times 10^3$ (25 points) CPU hours.

Consider first the setting with largest sample sizes: $n = 5$ for the high-fidelity GP and $n = 25$ for remaining models. Figure 8 (left) shows the predicted functions for each emulator, along with the high-fidelity test set. We see that both the M$^2$GP and the high-fidelity GP perform quite well in terms of predicting the high-fidelity function $\phi(x_1)$. Despite the similar predictive performance, the key difference between these methods (as we see later)
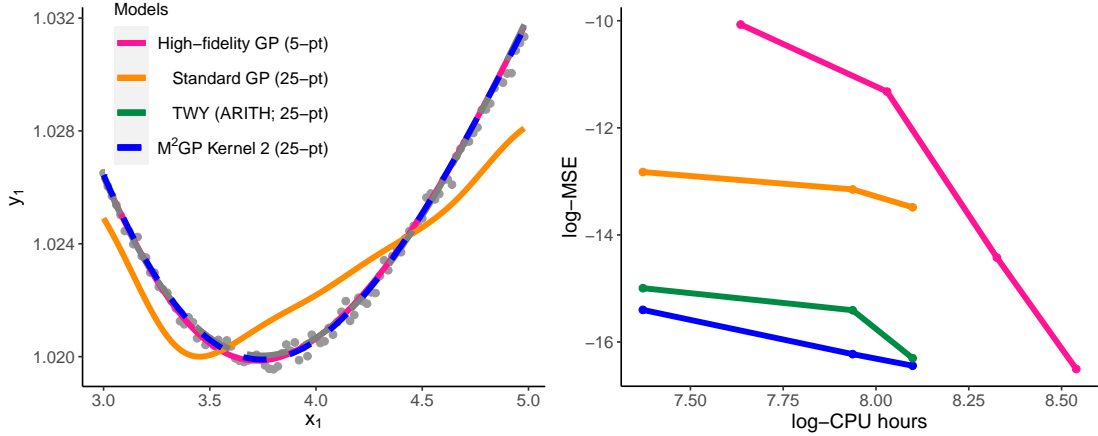
31

Figure 8: (Left) Predictions of $\phi(x_1)$ for the compared methods, along with the high-fidelity test data (grey). (Right) Plot of log-MSE vs. log-CPU hours required for simulating the training data for each emulator model.

is the cost of generating the required training data: this is much more expensive for the high-fidelity GP compared to the $M^2GP$, which affirms the value of multi-fidelity modeling. The TWY (ARITH) model provides slightly worse predictions, which is not surprising since it aggregates the complex multi-stage framework into a single fidelity parameter. The standard GP yields the worst performance of the compared methods.

Consider next the comparison of the predictive performance of the emulators given a computational budget for training data generation. Figure 8 (right) plots the log-MSEs of the considered models and their corresponding costs (in log-CPU hours) for simulating the training data. We see that, at a given computational budget, the $M^2GP$ yields the best predictive performance out of all the methods. This suggests that by integrating information on the underlying multi-stage multi-fidelity framework within its kernel specification, the proposed model can provide *cost-efficient* and accurate emulation of expensive simulators given a tight computational budget.

# 7    Conclusion

In this paper we presented a new emulator model, called the M$^2$GP, that tackles the challenge of surrogate modeling for multi-stage multi-fidelity simulators, whose precision is controlled by multiple fidelity parameters. Such simulators are often encountered in complex physical systems (including our motivating application in high-energy nuclear physics), but there has been little work in constructing cost-efficient emulators which leverage this structure for predictive modeling. The M$^2$GP makes use of novel non-stationary covariance functions, which embed numerical convergence information on the underlying multi-stage simulator within its kernel specification. This infusion of prior information allows for effective surrogate modeling of complex simulators, even with limited training data. We demonstrate the effectiveness of the M$^2$GP model in a suite of simulation experiments and in two applications, the first on emulating cantilever beam deflection, and the second on emulating the quark-gluon plasma in high-energy physics.

With these encouraging results, there are many avenues for future work. Given the promise of multi-fidelity modeling, one crucial direction for maximizing predictive power given a tight computational budget is experimental design. For example, given a budget of $10^6$ CPU hours for a project, an experimenter would wish to know if a better predictive model can be trained with a few carefully-chosen higher-fidelity runs, or with more lower-fidelity runs. Tackling this design problem for the current multi-stage multi-fidelity framework can greatly increase the applicability of the M$^2$GP in applications. We also aim to extend the M$^2$GP for a broader range of multi-fidelity applications, where simulator fidelity is more complex and cannot be well-captured by several continuous fidelity parameters.

# References

Brenner, S. C. & Scott, L. R. (2008), *The Mathematical Theory of Finite Element Methods*, Vol. 3, Springer.

Byrd, R. H., Lu, P., Nocedal, J. & Zhu, C. (1995), 'A limited memory algorithm for bound constrained optimization', *SIAM Journal on Scientific Computing* **16**(5), 1190–1208.

Cao, S., Chen, Y., Coleman, J., Mulligan, J., Jacobs, P., Soltz, R., Angerami, A., Arora, R., Bass, S., Cunqueiro, L. et al. (2021), 'Determining the jet transport coefficient $\hat{q}$ from inclusive hadron suppression measurements using Bayesian parameter estimation', *Physical Review C* **104**(2), 024905.

Chakraborty, A., Gopalakrishnan, S. & Reddy, J. N. (2003), 'A new beam finite element for the analysis of functionally graded materials', *International Journal of Mechanical Sciences* **45**(3), 519–539.

Chen, J., Mak, S., Joseph, V. R. & Zhang, C. (2021), 'Function-on-function kriging, with applications to three-dimensional printing of aortic tissues', *Technometrics* **63**(3), 384–395.

Chivers, C. (2012), *MHadaptive: General Markov Chain Monte Carlo for Bayesian Inference Using Adaptive Metropolis-Hastings Sampling*. R package version 1.1-8.
**URL:** *https://cran.r-project.org/package=MHadaptive*

Cox, D. D., Park, J.-S. & Singer, C. E. (2001), 'A statistical method for tuning a computer code to a data base', *Computational Statistics & Data Analysis* **37**(1), 77–92.

Cui, W. C. & Wisnom, M. R. (1992), 'Contact finite element analysis of three-and four-point short-beam bending of unidirectional composites', *Composites Science and Technology* **45**(4), 323–334.

Currin, C., Mitchell, T., Morris, M. & Ylvisaker, D. (1991), 'Bayesian prediction of deterministic functions, with applications to the design and analysis of computer experiments', *Journal of the American Statistical Association* **86**(416), 953–963.

De Souza, R. D., Koide, T. & Kodama, T. (2016), 'Hydrodynamic approaches in relativistic heavy ion reactions', *Progress in Particle and Nuclear Physics* **86**, 35–85.

DiazDelaO, F. A. & Adhikari, S. (2012), 'Bayesian assimilation of multi-fidelity finite element models', *Computers & Structures* **92**, 206–215.

Durrett, R. (2019), *Probability: Theory and Examples*, Vol. 49, Cambridge University Press.

Everett, D., Ke, W., Paquet, J. F., Vujanovic, G., Bass, S. A., Du, L., Gale, C., Heffernan, M., Heinz, U., Liyanage, D. et al. (2021), 'Multisystem Bayesian constraints on the transport coefficients of QCD matter', *Physical Review C* **103**(5), 054904.

Forrester, A. I., Sóbester, A. & Keane, A. J. (2007), 'Multi-fidelity optimization via surrogate modelling', *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* **463**(2088), 3251–3269.

Fricker, T. E., Oakley, J. E. & Urban, N. M. (2013), 'Multivariate Gaussian process emulators with nonseparable covariance structures', *Technometrics* **55**(1), 47–56.

Gale, C., Jeon, S. & Schenke, B. (2013), 'Hydrodynamic modeling of heavy-ion collisions', *International Journal of Modern Physics A* **28**(11), 1340011.

Gelman, A., Carlin, J. B., Stern, H. S. & Rubin, D. B. (1995), *Bayesian Data Analysis*, Chapman and Hall/CRC.

Gul, E., Joseph, V. R., Yan, H. & Melkote, S. N. (2018), 'Uncertainty quantification of machining simulations using an *in situ* emulator', *Journal of Quality Technology* **50**(3), 253–261.

He, X., Tuo, R. & Wu, C. F. J. (2017), 'Optimization of multi-fidelity computer experiments via the EQIE criterion', *Technometrics* **59**(1), 58–68.

Heinz, U. & Snellings, R. (2013), 'Collective flow and viscosity in relativistic heavy-ion collisions', *Annual Review of Nuclear and Particle Science* **63**, 123–151.

Heyliger, P. R. & Reddy, J. N. (1988), 'A higher order beam finite element for bending and vibration problems', *Journal of Sound and Vibration* **126**(2), 309–326.

Ho, M.-F., Bird, S. & Shelton, C. R. (2022), 'Multifidelity emulation for the matter power spectrum using Gaussian processes', *Monthly Notices of the Royal Astronomical Society* **509**(2), 2551–2565.

Hundsdorfer, W. H., Verwer, J. G. & Hundsdorfer, W. (2003), *Numerical Solution of Time-Dependent Advection-Diffusion-Reaction Equations*, Vol. 33, Springer.

Ji, Y., Mak, S., Soeder, D., Paquet, J. & Bass, S. A. (2021), 'A graphical Gaussian process model for multi-fidelity emulation of expensive computer codes', *arXiv preprint arXiv:2108.00306* .

Jin, S. S., Kim, S. T. & Park, Y. H. (2021), 'Combining point and distributed strain sensor for complementary data-fusion: A multi-fidelity approach', *Mechanical Systems and Signal Processing* **157**, 107725.

Johnson, M. E., Moore, L. M. & Ylvisaker, D. (1990), 'Minimax and maximin distance designs', *Journal of Statistical Planning and Inference* **26**(2), 131–148.

Joseph, V. R. (2016), 'Space-filling designs for computer experiments: A review', *Quality Engineering* **28**(1), 28–35.

Joseph, V. R., Gul, E. & Ba, S. (2015), 'Maximum projection designs for computer experiments', *Biometrika* **102**(2), 371–380.

Kennedy, M. C. & O'Hagan, A. (2000), 'Predicting the output from a complex computer code when fast approximations are available', *Biometrika* **87**(1), 1–13.

Kou, J. & Zhang, W. (2019), 'Multi-fidelity modeling framework for nonlinear unsteady aerodynamics of airfoils', *Applied Mathematical Modelling* **76**, 832–855.

Liyanage, D., Ji, Y., Everett, D., Heffernan, M., Heinz, U., Mak, S. & Paquet, J.-F. (2022), 'Efficient emulation of relativistic heavy ion collisions with transfer learning', *Physical Review C* **105**(3), 034910.

Mak, S. & Joseph, V. R. (2018), 'Minimax and minimax projection designs using clustering', *Journal of Computational and Graphical Statistics* **27**(1), 166–178.

Mak, S., Sung, C. L., Wang, X., Yeh, S. T., Chang, Y. H., Joseph, V. R., Yang, V. & Wu, C. F. J. (2018), 'An efficient surrogate model for emulation and physics extraction of large eddy simulations', *Journal of the American Statistical Association* **113**(524), 1443–1456.

Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H. & Teller, E. (1953), 'Equation of state calculations by fast computing machines', *The Journal of Chemical Physics* **21**(6), 1087–1092.

More, S. T. & Bindu, R. S. (2015), 'Effect of mesh size on finite element analysis of plate structure', *International Journal of Engineering Science and Innovative Technology* **4**(3), 181–185.

Morris, M. D. & Mitchell, T. J. (1995), 'Exploratory designs for computational experiments', *Journal of Statistical Planning and Inference* **43**(3), 381–402.

Ngo, D. & Scordelis, A. C. (1967), 'Finite element analysis of reinforced concrete beams', *Journal of the American Concrete Institute* **64**(3), 152–163.

Nocedal, J. & Wright, S. J. (1999), *Numerical Optimization*, Springer.

Park, S. J., Earmme, Y. Y. & Song, J. H. (1997), 'Determination of the most appropriate mesh size for a 2-d finite element analysis of fatigue crack closure behaviour', *Fatigue & Fracture of Engineering Materials & Structures* **20**(4), 533–545.

Perdikaris, P., Raissi, M., Damianou, A., Lawrence, N. D. & 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.

Perdikaris, P., Venturi, D., Royset, J. O. & Karniadakis, G. E. (2015), 'Multi-fidelity modelling via recursive co-kriging and Gaussian–Markov random fields', *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* **471**(2179), 20150018.

Qian, P. Z. & Wu, C. F. J. (2008), 'Bayesian hierarchical modeling for integrating low-accuracy and high-accuracy experiments', *Technometrics* **50**(2), 192–204.

Qian, Z., Seepersad, C. C., Joseph, V. R., Allen, J. K. & Wu, C. F. J. (2006), 'Building surrogate models based on detailed and approximate simulations', *Journal of Mechanical Design* **128**(4), 668–677.

Reese, C. S., Wilson, A. G., Hamada, M., Martz, H. F. & Ryan, K. J. (2004), 'Integrated analysis of computer and physical experiments', *Technometrics* **46**(2), 153–164.

Santner, T. J., Williams, B. J., Notz, W. I. & Williams, B. J. (2003), *The Design and Analysis of Computer Experiments*, Vol. 1, Springer.

Shi, R., Liu, L., Long, T., Wu, Y. & Wang, G. G. (2020), 'Multi-fidelity modeling and adaptive co-kriging-based optimization for all-electric geostationary orbit satellite systems', *Journal of Mechanical Design* **142**(2).

Smith, M. (2009), *ABAQUS/Standard User's Manual, Version 6.9*, Dassault Systèmes Simulia Corp, United States.

Stein, A. & Corsten, L. (1991), 'Universal kriging and cokriging as a regression procedure', *Biometrics* **47**(2), 575–587.

Sun, F., Gramacy, R. B., Haaland, B., Lu, S. & Hwang, Y. (2019), 'Synthesizing simulation and field data of solar irradiance', *Statistical Analysis and Data Mining: The ASA Data Science Journal* **12**(4), 311–324.

Templeton, J. A., Blaylock, M. L., Domino, S. P., Hewson, J. C., Kumar, P. R., Ling, J., Najm, H. N., Ruiz, A., Safta, C., Sargsyan, K. et al. (2015), Calibration and forward uncertainty propagation for large-eddy simulations of engineering flows, Technical report, Sandia National Lab., Livermore, CA.

Tuo, R., Wu, C. F. J. & Yu, D. (2014), 'Surrogate modeling of computer experiments with different mesh densities', *Technometrics* **56**(3), 372–380.

Vanden-Eijnden, E. (2003), 'Numerical techniques for multi-scale dynamical systems with stochastic effects', *Communications in Mathematical Sciences* **1**(2), 385–391.

Williams, C. K. & Rasmussen, C. E. (2006), *Gaussian Processes for Machine Learning*, Vol. 2, MIT Press Cambridge, MA.

Wu, C. F. J. & Hamada, M. S. (2009), *Experiments: Planning, Analysis, and Optimization*, John Wiley & Sons.

Xiong, S., Qian, P. Z. & Wu, C. F. J. (2013), 'Sequential design and analysis of high-accuracy and low-accuracy computer codes', *Technometrics* **55**(1), 37–46.