

# A VARIATIONAL APPROACH TO PHYSICS INFORMED NEURAL NETWORK FOR STOCHASTIC PARTIAL DIFFERENTIAL EQUATIONS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Physics-Informed Neural Networks (PINNs) refer to a deep learning approach to represent the spatial and temporal characteristics of a distributed physical phenomenon, such as thermal fields, using Neural Networks (NNs). The loss function used to train PINNs relies on a term that penalises any violation of the Partial Differential Equation (PDE) that governs the distributed phenomena, in addition to a least-squares-error penalisation of any available observations in the relevant domain. PINNs have been shown to be successful in approximating the solutions to PDEs and have proven effective even in the low data regime. A critical shortcoming is the lack of a systematic treatment of the uncertainty of the approximation. State-of-the-art approaches rely on Bayesian NNs, though these are computationally heavy, both during training and at inference, nor does the associated training procedure derive from first principles. To remedy these limitations, we propose Variational Inference PINNs (VI-PINNs). Our approach derives from first principles: the uncertainty intrinsic to the distributed phenomena is explained by adopting Stochastic PDEs, while relying on standard measurement uncertainty to explain the observational uncertainty. This leads to the formulation of a posterior probability for the distributed phenomenon. Drawing parallels with Bayesian inference in finite spaces and relying on VI techniques to circumvent the otherwise intractable posterior. We derive a training objective that allows us to train two NNs, representing the mean and covariance of the approximation, respectively. The solution may be interpreted as a Bayesian belief about the true distributed phenomenon. Importantly, in the limit, the original PINN framework is recovered. We compare our approach with Bayesian PINNs (B-PINNs). Our results suggest that VI-PINNs are easier to implement, have a lower training time, and yields results that better align with reality, especially when extrapolating outside the measurement range.

**Keywords:** Neural Network, Variational Inference, Steady State Heat Equation, Physics-Informed Neural Network, Partial Differential Equations, Gaussian Process

## 1 INTRODUCTION

A Physics-Informed Neural Network (PINN) is a new class of Neural Network (NN) introduced by Raissi et al. (2017). The main idea of a PINN is to use the strengths of a NN, such as approximating (complex) models [Markidis (2021), Bhatnagar et al. (2024), Fowler et al. (2024)], modelling unknown parts of a model [Markidis (2021)], or pattern recognition in data sets [Mah & Chakravarthy (1992)]. That is done in combination with the underlying physics of a system. The underlying physics is introduced through the use of a Partial Differential Equations (PDEs) or Ordinary Differential Equations (ODEs). This allows for a smaller dataset of measurements, compared to a traditional NN. It can be used to solve forward problems, such as approximating/ emulating solutions to PDEs and ODEs, or inverse problems, such as estimating model parameters from limited data.

PINNs have been shown to yield good results on a broad range of applications [Raissi et al. (2017), Markidis (2021)], but they also have their limitations. PINNs has no uncertainty quantification built into it. This makes PINNs sensitive to the noise in the training data. In Yang et al. (2021),

a Bayesian PINN (B-PINN) is proposed to solve this problem. The uncertainty is quantified using sample-based techniques such as Hamiltonian Monte Carlo (HMC) [Cobb & Jalaian (2021)] or No-U-Turn Sampler (NUTS) [Hoffman & Gelman (2011)]. These are computation-heavy and result in a long training time for each iteration step.

To overcome this problem, we propose a Variational Inference Physics Informed Neural Network (VI-PINN). This topology quantifies the uncertainty in an algebraic way instead of using sampling. This paper is structured as follows: In section 2, the background for PINN and B-PINN is given. In section 3, the VI-PINN is constructed using Variational Inference and Bayesian inference. This is first done for the general finite-dimensional case, followed by an ad hoc generalisation. In section 4, the theoretical framework is implemented in a 1D problem, where it will be compared with the results of a B-PINN. In section 5, conclusions are presented.

## 2 BACKGROUND

### 2.1 PHYSICS INFORMED NEURAL NETWORK

The dynamics of a system can be represented using PDEs, which has the general form

$$\mathcal{A}_{\mathbf{x},t}[\mathbf{u}](\mathbf{x}, t) = \partial_t \mathbf{u}(\mathbf{x}, t) + \mathcal{A}_{\mathbf{x}}[\mathbf{u}](\mathbf{x}, t) = 0 \quad (1)$$

$\mathbf{x} \in \mathbb{R}^n$  is a spatial vector,  $t \in \mathbb{R}$  the time,  $\partial_t \mathbf{u}(\mathbf{x}, t)$  is the partial derivative of  $\mathbf{u}(\mathbf{x}, t)$  with respect to  $t$  and  $\mathcal{A}_{\mathbf{x}}$  is some spatial linear operator. The objective is to find the explicit formulation  $\mathbf{u}(\mathbf{x}, t)$  that satisfies equation (1), so that  $\mathbf{u}(\mathbf{x}, t)$  describes the dynamics in the spatio-temporal domain.

In a PINN, the unknown function  $\mathbf{u}$  is approximated by the NN, with as output  $\mathbf{u}_{\theta}$  and NN-parameters  $\theta$ . The cost function is of the form

$$J_{cost}(\theta) = (1 - \lambda) \sum_{(\mathbf{x}_m, t_m, \mathbf{y}) \in \mathcal{D}_{DATA}} \frac{1}{2} \|\mathbf{y} - \mathbf{u}_{\theta}(\mathbf{x}_m, t_m)\|^2 + \lambda \sum_{(\mathbf{x}, t) \in \mathcal{D}_{PHYS}} \frac{1}{2} \|\mathcal{A}_{\mathbf{x},t}[\mathbf{u}_{\theta}](\mathbf{x}, t)\|^2 \quad (2)$$

$\mathbf{y}$  are measurements at  $(\mathbf{x}_m, t_m)$ ,  $\lambda$  is a tuning factor that determines the relative importance of each cost function,  $\mathcal{D}_{DATA}$  is the data set, and  $\mathcal{D}_{PHYS}$  represent an arbitrary set of space-time coordinates where the PDEs is enforced to the NN. In the literature, there is no clear way to determine the value of  $\lambda$ , nor the set of space-time coordinates.

### 2.2 BAYESIAN PHYSICS-INFORMED NEURAL NETWORK

A way to quantify the uncertainty of the prediction is by extending the PINN framework to a B-PINN [Yang et al. (2021)]. This can be done by following the Bayesian Neural Network (BNN) approach [Jospin et al. (2022)] and constructing a posterior for the NN-parameters (and optionally the PDEs parameters).

$$p(\theta | \mathcal{D}_{DATA}, \mathcal{D}_{PHYS}) \propto p(\mathcal{D}_{DATA} | \theta) p(\mathcal{D}_{PHYS} | \theta)$$

where

$$p(\mathcal{D}_{DATA} | \theta) = \prod_{(\mathbf{x}_m, t_m, \mathbf{y}) \in \mathcal{D}_{DATA}} \mathcal{N}(\mathbf{y} | \mathbf{u}_{\theta}(\mathbf{x}_m, t_m), \sigma_{\mathbf{y}}^2) \quad (3a)$$

$$p(\mathcal{D}_{PHYS} | \theta) = \prod_{(\mathbf{x}, t) \in \mathcal{D}_{PHYS}} \mathcal{N}(0 | \mathcal{A}_{\mathbf{x},t}[\mathbf{u}_{\theta}](\mathbf{x}, t), \sigma_{\mathbf{u}}^2) \quad (3b)$$

Since the posterior is intractable, it can be approximated using Markov Chain Monte Carlo sampling. Possible strategies to establish these samples are HMC or NUTS. In Figure 1, the B-PINN approach is schematically shown.

In the end, one obtains a weighted particle set  $\{(\mathbf{a}_i, \theta_i)\}$  that approximates  $p(\theta | \mathcal{D}_{DATA}, \mathcal{D}_{PHYS})$ . This particle set can be used to evaluate the output uncertainty of the distribution  $\mathbf{u}$ , given  $(\mathbf{x}, t)$ , according to the predicted distribution

$$p(\mathbf{u} | \mathbf{x}, t) = \int p(\theta | \mathcal{D}_{DATA}, \mathcal{D}_{PHYS}) \delta(\mathbf{u} - \mathbf{u}_{\theta}(\mathbf{x}, t)) d\theta \approx \sum_i \mathbf{a}_i \delta(\mathbf{u} - \mathbf{u}_{\theta_i}(\mathbf{x}, t))$$

and certain statistical measures in particular, e.g.

$$\mu_k(\mathbf{x}, t) = \int \mathbf{u}^k p(\mathbf{u} | \mathbf{x}, t) d\mathbf{u} \approx \sum_i \mathbf{a}_i \mathbf{u}_{\theta_i}(\mathbf{x}, t)^k \quad (4)$$

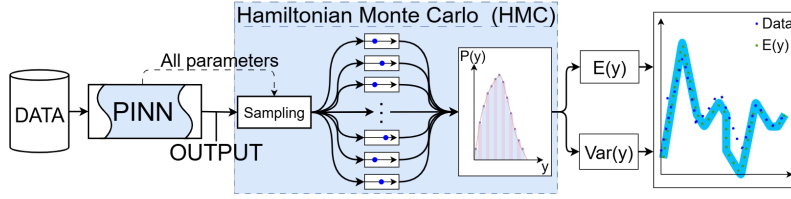


Figure 1: Schematic of a B-PINN using HMC

### 2.2.1 IMPLEMENTATION OF THE B-PINN

A B-PINN is capable of finding a good approximation, but it does so inefficiently. It takes  $N_{MC}$  samples of the uncertain parameters. This includes all weights and biases from the NN, and optionally the parameters of the PDEs too. All  $N_{par}$  uncertain parameters are assumed to have a form of uncertainty [Cobb & Jalaian (2021)]. This results in  $N_{MC} \cdot N_{par}$  samples being taken for every iteration step. A proper estimation of the output uncertainty requires  $N_{it}$  training iterations of  $\mathbf{u}_\theta$  to get a reliable result for the mean and covariance (i.e.  $k=1$  and  $k=2$  in equation (4), respectively).

When the complexity and/or noise level increases, the total required number of samples [ $N_{it} \cdot (N_{MC} \cdot N_{par})$ ] will increase quickly. We desire to derive a more efficient approach. Our idea is to approximate  $p(\mathbf{u})$  using a variational density and learn the density using Variational Inference (VI).

### 2.3 LINK WITH STOCHASTIC PARTIAL DIFFERENTIAL EQUATIONS

As one can note, Stochastic Partial Differential Equations (SPDEs) (equation (5)) are linked with equation (3b).

$$\mathcal{A}_{x,t}[\mathbf{u}](\mathbf{x}, t) = \partial_t \mathbf{u}(\mathbf{x}, t) + \mathcal{A}_x[\mathbf{u}](\mathbf{x}, t) = \omega(\mathbf{x}, t) \quad (5)$$

with  $\mathbf{u}$  is the solution of equation (5),  $\omega(\mathbf{x}, t)$  the stochastic force term.  $\omega$  can be modelled using a Gaussian Process (GP), with mean  $\mu_\omega$ , and kernel  $\sigma_\omega^2$ .

$$\omega \sim \mathcal{GP}(\mu_\omega, \sigma_\omega^2)$$

In case that  $\mathcal{A}_x$  is a linear operator, from Rasmussen & Williams (2006) it is known that the function  $\mathbf{u}$  is also a GP where the mean and covariance function are governed by the SPDEs.

$$\mathbf{u}(\mathbf{x}, t) \sim \mathcal{GP}(\mu_{\mathbf{u}}, \sigma_{\mathbf{u}}^2)$$

## 3 A VARIATIONAL APPROACH TO PHYSICS INFORMED NEURAL NETWORKS

### 3.1 BAYES' RULE

The posterior for the quantity  $\mathbf{u}$ , given the measurements  $\mathbf{y}$ , is given by *Bayes' rule* [Bolstad & Curran (2017)]. We assume to have access to the measurement model  $p(\mathbf{y}, \mathbf{u})$ , and prior  $p(\mathbf{u})$ .

$$p(\mathbf{u}|\mathbf{y}) = \frac{p(\mathbf{y}, \mathbf{u})}{p(\mathbf{y})} = \frac{p(\mathbf{y}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{y})} \quad (6)$$

#### 3.1.1 FINITE-DIMENSIONAL BAYESIAN INFERENCE

Assume that the quantity  $\mathbf{u}$  is actually a function, describe by equation (7) where  $\mathbf{A}$  is non-singular and  $\omega \sim \mathcal{GP}(\mu_\omega, \Sigma_{\omega\omega})$ .

$$\mathbf{A}\mathbf{u} = \omega \quad (7)$$

Due to the linearity of equation (7), it is known that [Sarkka & Hartikainen (2012), Lindgren et al. (2011)]

$$\mathbf{u} \sim p(\mathbf{u}) = \mathcal{GP}(\mu_{\mathbf{u}}, \Sigma_{\mathbf{u}\mathbf{u}})$$

where

$$\begin{aligned} \mu_{\mathbf{u}} &= \mathbb{E}[\mathbf{u}] = \mathbb{E}[\mathbf{A}^{-1}\omega] = \mathbf{A}^{-1}\mu_\omega \\ \Sigma_{\mathbf{u}\mathbf{u}} &= \mathbb{E}[(\mathbf{u} - \mu_{\mathbf{u}})(\mathbf{u} - \mu_{\mathbf{u}})^T] = \mathbb{E}[\mathbf{A}^{-1}(\omega - \mu_\omega)(\omega - \mu_\omega)^T \mathbf{A}^{-T}] = \mathbf{A}^{-1}\Sigma_{\omega\omega}\mathbf{A}^{-T} \end{aligned} \quad (8)$$

162 Additionally, we assume that a measurement is taken

$$163 \mathbf{y} \sim \mathcal{N}(\mathbf{B}\mathbf{u}, \Sigma_{\mathbf{y}\mathbf{y}})$$

164 such that

$$165 \begin{bmatrix} \mathbf{u} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{GP} \left( \begin{bmatrix} \mu_{\mathbf{u}} \\ \mathbf{B}\mu_{\mathbf{u}} \end{bmatrix}, \begin{bmatrix} \Sigma_{\mathbf{u}\mathbf{u}} & \Sigma_{\mathbf{u}\mathbf{u}}\mathbf{B}^T \\ \mathbf{B}\Sigma_{\mathbf{u}\mathbf{u}} & \Sigma_{\mathbf{y}\mathbf{y}} + \mathbf{B}\Sigma_{\mathbf{u}\mathbf{u}}\mathbf{B}^T \end{bmatrix} \right) \quad (9)$$

166 To conclude, the conditional probability for the means and variance is described by

$$167 \begin{aligned} 168 \mu_{\mathbf{u}|\mathbf{y}} &= \mu_{\mathbf{u}} + \Sigma_{\mathbf{u}\mathbf{u}}\mathbf{B}^T(\Sigma_{\mathbf{y}\mathbf{y}} + \mathbf{B}\Sigma_{\mathbf{u}\mathbf{u}}\mathbf{B}^T)^{-1}(\mathbf{y} - \mathbf{B}\mu_{\mathbf{u}}) \\ 169 \Sigma_{\mathbf{u}\mathbf{u}|\mathbf{y}} &= \Sigma_{\mathbf{u}\mathbf{u}} - \Sigma_{\mathbf{u}\mathbf{u}}\mathbf{B}^T(\Sigma_{\mathbf{y}\mathbf{y}} + \mathbf{B}\Sigma_{\mathbf{u}\mathbf{u}}\mathbf{B}^T)^{-1}\mathbf{B}\Sigma_{\mathbf{u}\mathbf{u}} \end{aligned} \quad (10)$$

### 173 3.2 VARIATIONAL INFERENCE FOR SPDE

174 *Bayes' rule* is elegant, but quite often not applicable in practice due to the unknown measurement distribution  $p(\mathbf{y})$ . A popular technique that is used to approximate the posterior is VI [Ganguly & Earp (2021)]. The idea is to fit a variational density  $q$ , which is a member of some density family  $\mathcal{Q}$ , by solving an optimisation problem. By minimising the Variational Free Energy (VFE) [Wu et al. (2019)],  $q$  can be found.

$$180 q^*(\mathbf{u}) = \arg \min_{q \in \mathcal{Q}} \mathbb{E}_{q(\mathbf{u})} \left[ \log \left( \frac{q(\mathbf{u})}{p(\mathbf{y}, \mathbf{u})} \right) \right] \quad (11)$$

181 The quantity  $\mathbf{u}$  is actually a function, where the entries of the vector  $\mathbf{u}$  can be seen as different function evolutions. The objective is now to find a variational density representing the probability of the function  $\mathbf{u}$ , that satisfies the SPDEs, and, additionally, conditions the likelihood on the measurements,  $\mathbf{y}$ .

182 The discussion will be limited to static SPDEs, though we note that the approach generalises trivially to spatio-temporal SPDEs.

#### 183 3.2.1 FINITE-DIMENSIONAL VARIATIONAL INFERENCE

184 We start reviewing a general finite-dimensional case and comparing it to the results from Bayesian Inference (section 3.1.1), with the objective to obtain the same results as equation (8).

185 Using a variational approximation, following Opper (2019), this can be found by introducing a variational density

$$186 \mathbf{u} \sim q(\mathbf{u}) = \mathcal{GP}(\mu_q, \Sigma_{qq})$$

187 with the free energy is given by

$$188 \mathcal{L}_{VI}[\mathbf{u}] = \mathbb{E}_{q(\mathbf{u})} \left[ \log \left( \frac{q(\mathbf{u})}{p(\mathbf{u})} \right) \right] \quad (12)$$

189 where

$$190 \begin{aligned} 191 q(\mathbf{u}) &= |2\pi\Sigma_{qq}|^{-\frac{1}{2}} \exp \left( -\frac{1}{2}(\mathbf{u} - \mu_q)^T \Sigma_{qq}^{-1}(\mathbf{u} - \mu_q) \right) \\ 192 p(\boldsymbol{\omega}|\mathbf{u}) \equiv p(\mathbf{u}) &= |2\pi\Sigma_{\boldsymbol{\omega}\boldsymbol{\omega}}|^{-\frac{1}{2}} \exp \left( -\frac{1}{2}(\mathbf{A}\mathbf{u} - \mu_{\boldsymbol{\omega}})^T \Sigma_{\boldsymbol{\omega}\boldsymbol{\omega}}^{-1}(\mathbf{A}\mathbf{u} - \mu_{\boldsymbol{\omega}}) \right) \end{aligned} \quad (13)$$

193 The VFE can be evaluated to find equation (14). The stepwise calculation can be found in Appendix B.1.

$$194 \mathcal{L}_{VI}[q] \propto -\frac{1}{2} \log(|\Sigma_{qq}|) + \frac{1}{2} \text{tr}(\Sigma_{\boldsymbol{\omega}\boldsymbol{\omega}}^{-1} \mathbf{A} \Sigma_{qq} \mathbf{A}^T) + \frac{1}{2} (\mathbf{A}\mu_q - \mu_{\boldsymbol{\omega}})^T \Sigma_{\boldsymbol{\omega}\boldsymbol{\omega}}^{-1} (\mathbf{A}\mu_q - \mu_{\boldsymbol{\omega}}) \quad (14)$$

195 To find the minimum of  $\mathcal{L}_{VI}[q]$ , the first partial derivative with respect to  $\mu_q$  and  $\Sigma_{qq}$  are set equal to 0. Using Petersen & Pedersen (2012), this results in equation (15). One verifies that equation (8) and equation (15) give the same result.

$$196 \mu_q = \mathbf{A}^{-1} \mu_{\boldsymbol{\omega}} \quad \Sigma_{qq} = \mathbf{A}^{-1} \Sigma_{\boldsymbol{\omega}\boldsymbol{\omega}} \mathbf{A}^{-T} \quad (15)$$

When adding additional measurements, equation (12) can be extended to

$$\mathcal{L}_{VI}^y[\mathbf{u}] = \mathbb{E}_{q(\mathbf{u})} \left[ \log \left( \frac{q(\mathbf{u})}{p(\mathbf{u}, \mathbf{y})} \right) \right] \quad (16)$$

and results in

$$\begin{aligned} \mathcal{L}_{VI}^y[q] \propto & \underbrace{\frac{1}{2}(\mathbf{A}\mu_q - \mu_\omega)^T \Sigma_{\omega\omega}^{-1}(\mathbf{A}\mu_q - \mu_\omega)}_{(17.1)} + \underbrace{\frac{1}{2}tr(\Sigma_{\omega\omega}^{-1}\mathbf{A}\Sigma_{qq}\mathbf{A}^T)}_{(17.2)} \\ & + \underbrace{\frac{1}{2}(\mathbf{B}\mu_q - \mathbf{y})^T \Sigma_{yy}^{-1}(\mathbf{B}\mu_q - \mathbf{y})}_{(17.3)} + \underbrace{\frac{1}{2}tr(\Sigma_{yy}^{-1}\mathbf{B}\Sigma_{qq}\mathbf{B}^T)}_{(17.4)} - \underbrace{\frac{1}{2}\log|\Sigma_{qq}|}_{(17.5)} \end{aligned} \quad (17)$$

The stepwise calculation can be found in Appendix B.2. One can verify that the minimum of equation (17), found by setting the partial derivatives equal to 0, gives the same results as equation (9).

Once we have access to the mean and covariance, in fact, the job is done, and one can come back to regular inference with GP. There exist standard techniques to find the mean and covariance for given SPDEs, however, all of them depend on linear function decompositions (e.g. Eigenfunctions, Finite Element Method (FEM), Finite Difference Method (FDM), ...) of the function  $\mathbf{u}$ , and shift all uncertainty to the coefficients. Rather, we would want to draw a connection with the PINN paradigm, specifically by approximating  $\mu_{\mathbf{u}}$  and  $\sigma_{\mathbf{u}}^2$ .

The ambition is to approach an SPDEs similar to how a PINN approaches a PDEs.

### 3.2.2 AD-HOC GENERALISATION TO INFINITE-DIMENSIONAL CASE

The next step is to establish an ad hoc generalisation to the infinite-dimensional case. This is done using the two analogies between: a vector with infinite vector entries and function evaluations at arbitrary arguments, and matrices and linear operators.

Consider the static SPDEs in equation (18), where  $\mathcal{A}_{\mathbf{x}}$  is some linear operator and  $\omega \sim \mathcal{GP}(\mu_\omega, \sigma_\omega^2)$ .

$$\mathcal{A}_{\mathbf{x}}[\mathbf{u}] = \omega \quad (18)$$

Due to the linearity, we have  $\mathbf{u} \sim \mathcal{GP}(\mu_{\mathbf{u}}, \sigma_{\mathbf{u}}^2)$ . Further assume that we take point measurements  $\mathbf{y}$ , at spatial coordinates  $\mathbf{x}_m$ , contained in some data set  $\mathcal{D}_{DATA}$ .

$$\mathbf{y} \sim \mathcal{N}(\mathcal{B}_{\mathbf{x}}[\mathbf{u}], \sigma_{\mathbf{y}}^2) \quad \mathcal{D}_{DATA} = \{(\mathbf{x}_m, \mathbf{y}_i) | i = 1 \dots N_m\}$$

where:

$$\mathcal{B}_{\mathbf{x}}[\mathbf{u}] = \mathbf{u}(\mathbf{x}_m)$$

or

$$\mathcal{B}_{\mathbf{x}}[\mathbf{u}] = \int \mathbf{u}(\mathbf{x}') \delta(\mathbf{x} - \mathbf{x}') d\mathbf{x}'$$

We already know that

$$\mathbf{u} \sim \mathcal{GP}(\mu_{\mathbf{u}}, \sigma_{\mathbf{u}}^2)$$

Using Opper (2019), this allows us to write the probability of a realisation of  $\mathbf{u}$  as

$$-\log p(\mathbf{u}) \propto \exp\left(-\frac{1}{2} \iint [\mathbf{u}(\mathbf{x}) - \mu_{\mathbf{u}}(\mathbf{x})]^T \sigma_{\mathbf{u}}^{-2}(\mathbf{x}, \mathbf{x}') [\mathbf{u}(\mathbf{x}) - \mu_{\mathbf{u}}(\mathbf{x})] d\mathbf{x} d\mathbf{x}'\right) \quad (19)$$

and the probability of the measurement data is written as

$$p(\mathcal{D}_{DATA} | \mathbf{u}) = \prod_{(\mathbf{x}_m, \mathbf{y}) \in \mathcal{D}_{DATA}} \left( \frac{1}{\sqrt{2\pi\sigma_{\mathbf{y}}^2}} \exp\left(-\frac{1}{2\sigma_{\mathbf{y}}^2} \|\mathbf{y} - \mathcal{B}_{\mathbf{x}}[\mathbf{u}]\|^2\right) \right) \quad (20)$$

Next, the ad-hoc generalisation of equation (17) is constructed.

$$(17.1) \Rightarrow \frac{1}{2} \iint_{\mathcal{D}_{PHYS}^2} \mathcal{A}_x[\mu_q]^T \sigma_\omega^{-2}(\mathbf{x}, \mathbf{x}') \mathcal{A}_x[\mu_q] d\mathbf{x} d\mathbf{x}' \quad (21)$$

$$(17.2) \Rightarrow \frac{1}{2} \iint_{\mathcal{D}_{PHYS}^2} tr(\sigma_\omega^{-2} \mathcal{A}_x[\sigma_q^2(\mathbf{x}, \mathbf{x}') \mathcal{A}_{x'}^T]) d\mathbf{x} d\mathbf{x}' \quad (22)$$

$$(17.3) \Rightarrow \frac{1}{2} \sum_{(\mathbf{x}_m, \mathbf{y}) \in \mathcal{D}_{DATA}} \sigma_y^{-2} \|\mathbf{y} - \mathcal{B}_x[\mu_q]\|^2 \quad (23)$$

$$(17.4) \Rightarrow \frac{1}{2} \sum_{(\mathbf{x}_m, \mathbf{y}) \in \mathcal{D}_{DATA}} tr(\sigma_y^{-2} \mathcal{B}_x[\sigma_q^2(\mathbf{x}_m, \mathbf{x}'_m)] \mathcal{B}_{x'}^T) \quad (24)$$

For (17.5), an ad-hoc solution is a great challenge. We proposed that it can be seen as a regularisation term for  $\sigma_q^2$ . If  $\sigma_q^2$  were a matrix, then would it be  $\log(\det[\sigma_q^2])$ . Here  $\sigma_q^2$  is a covariance kernel, the generalised determinant,  $\det[\cdot]$ , it is then basically the product of the eigenvalues of the kernel. In practice, only a finite number of query points are being evaluated. The determinant of the associated covariance matrix can be approximated by

$$(17.5) \Rightarrow \det[\sigma_q^2] \approx \det(\Sigma_{qq}^{PHYS}) \quad (25)$$

where  $[\Sigma_{qq}^{PHYS}]_{ij} = \sigma_q^2(\mathbf{x}_i, \mathbf{x}_j)$ ;  $\mathbf{x}_i, \mathbf{x}_j \in \mathcal{D}_{PHYS}$

The last obstacles in the derivation are the double integrals. Those are impossible to evaluate for arbitrary networks. Hence, the double integral can be approximated by drawing random samples in the spatial dimension.

$$\iint_{\mathcal{D}^2} f(\mathbf{x}, \mathbf{x}') d\mathbf{x} d\mathbf{x}' \approx \Delta \mathbf{x}^2 \sum_{(\mathbf{x}, \mathbf{x}') \in \mathcal{D}^2} f(\mathbf{x}, \mathbf{x}') \quad (26)$$

Filling equation (21-26) into equation (17), the approximated total free energy can be calculated

$$\begin{aligned} \mathcal{L}_{VI}^y[\mu_q, \sigma_q^2] &\approx \underbrace{\frac{1}{2} \Delta \mathbf{x}^2 \sum_{(\mathbf{x}, \mathbf{x}') \in \mathcal{D}_{PHYS}^2} \mathcal{A}_x[\mu_q]^T \sigma_\omega^{-2}(\mathbf{x}, \mathbf{x}') \mathcal{A}_{x'}[\mu_q]}_{J_1[\mu_q]} + \underbrace{\frac{1}{2} \Delta \mathbf{x}^2 \sum_{(\mathbf{x}, \mathbf{x}') \in \mathcal{D}_{PHYS}^2} tr(\sigma_\omega^{-2} \mathcal{A}_x[\sigma_q^2(\mathbf{x}, \mathbf{x}') \mathcal{A}_{x'}^T])}_{J_2[\sigma_q^2]} \\ &+ \underbrace{\frac{1}{2} \sum_{(\mathbf{x}_m, \mathbf{y}) \in \mathcal{D}_{DATA}} tr(\sigma_y^{-2} \mathcal{B}_x[\sigma_q^2(\mathbf{x}_m, \mathbf{x}'_m)] \mathcal{B}_{x'}^T)}_{J_3[\sigma_q^2]} + \underbrace{\frac{1}{2} \sum_{(\mathbf{x}_m, \mathbf{y}) \in \mathcal{D}_{DATA}} \sigma_y^{-2} \|\mathbf{y} - \mathcal{B}_x[\mu_q]\|^2}_{J_4[\mu_q]} - \underbrace{\frac{1}{2} \log |\det[\Sigma_{qq}^{PHYS}]|}_{J_5[\sigma_q^2]} \end{aligned} \quad (27)$$

### 3.3 RELATION WITH PHYSICS INFORMED NEURAL NETWORK

Equation (27) can be split into two parts. A part that is similar to a standard PINN cost function (equation (2) w.r.t.  $\mu_q$ ).

$$J_1[\mu_q] + J_4[\mu_q] = \frac{1}{2} \Delta \mathbf{x}^2 \sum_{(\mathbf{x}, \mathbf{x}') \in \mathcal{D}_{PHYS}^2} \mathcal{A}_x[\mu_q](\mathbf{x})^T \sigma_\omega^{-2}(\mathbf{x}, \mathbf{x}') \mathcal{A}_{x'}[\mu_q](\mathbf{x}') + \frac{1}{2} \sum_{(\mathbf{x}_m, \mathbf{y}) \in \mathcal{D}_{DATA}} \sigma_y^{-2} \|\mathbf{y} - \mathcal{B}_x[\mu_q]\|^2$$

and a second part that expresses the uncertainty (w.r.t.  $\sigma_q^2$ ).

$$\begin{aligned} J_2[\sigma_q^2] - J_5[\sigma_q^2] + J_3[\sigma_q^2] &= \frac{1}{2} \Delta \mathbf{x}^2 \sum_{(\mathbf{x}, \mathbf{x}') \in \mathcal{D}_{PHYS}^2} tr(\sigma_\omega^{-2} \mathcal{A}_x[\sigma_q^2(\mathbf{x}_m, \mathbf{x}'_m)] \mathcal{A}_{x'}^T) - \frac{1}{2} \log |\det[\Sigma_{qq}^{PHYS}]| \\ &+ \frac{1}{2} \sum_{(\mathbf{x}_m, \mathbf{y}) \in \mathcal{D}_{DATA}} tr(\sigma_y^{-2} \mathcal{B}_x[\sigma_q^2(\mathbf{x}_m, \mathbf{x}'_m)] \mathcal{B}_x^T) \end{aligned} \quad (28)$$

As one can see, there is no explicit tuning factor,  $\lambda$ . The tuning results from  $\sigma_y^2$ ,  $\sigma_\omega^2$  and  $\Delta x$ , i.e. the spatial coarseness.

## 4 RESULTS

### 4.1 PROBLEM FORMULATION

To validate our proposed framework, a 1D steady state Disturbed Thermal Field (DTF) is examined, which can be described by means of the following SPDEs [Fourier (1878)] and Boundary Condition (BC):

$$k \frac{\partial^2 \mathbf{u}(\mathbf{x})}{\partial \mathbf{x}^2} = \omega(\mathbf{x}) \quad \text{BC : } \mathbf{u}(0) = \mathbf{u}(1) = 1 \quad (29)$$

with  $\mathbf{x} \in [0, 1.0]$ , the spatial coordinates,  $k > 0$ , the conductivity constant,  $\omega \sim \mathcal{GP}(\mu_\omega, \sigma_\omega^2)$ , the stochastic heat flux term and  $\mathbf{u} \sim \mathcal{GP}(\mu_u, \sigma_u^2)$ , the stochastic temperate field. Take a 1D steel beam with a length of 1 m with homogeneous material properties. Over the length of the beam some torches are placed, as shown in Figure 2. Normalizing equation (29), such that  $k = 1$ , results in a DTF with  $\mu_\omega = (6\pi)^2 \sin(6\pi \cdot \mathbf{x})$ , and variance,  $\sigma_\omega^2 = 4, 0 \cdot \delta(\mathbf{x} - \mathbf{x}')$ .

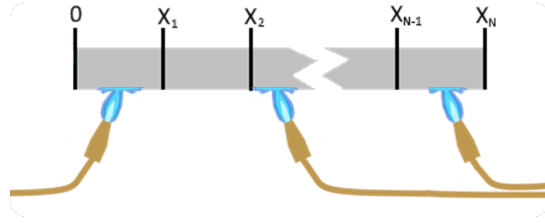


Figure 2: 1D steel beam with distributed heat sources

For the training of the different NNs, 25  $x_m$ -points are chosen evenly spread over  $[0, 0.6]$ . For the testing/validation, 200  $x$ -points are chosen evenly spread over  $[0, 1.0]$ . This allows us to compare the inter-/and extrapolation capabilities of the B-PINN and VI-PINN. The training points have some noise  $\epsilon_y \sim \mathcal{GP}(0, 1.5^2)$ , added to emulate measurement noise.

It was found empirically that rescaling the input  $x$ , and measurement data  $y$ , yielded better results and reduced computation time.

$$\mathcal{D}_{DATA} \in \mathbb{R}^n \mapsto \mathcal{D}'_{DATA} \in [0, 1]^n$$

### 4.2 SOLVE WITH A B-PINN USING NUTS

The implementation of the HMC technique is performed using the NUTS [Hoffman & Gelman (2011)] implementation in Python using NumPyro [Phan et al. (2019), Bingham et al. (2019)]. This algorithm is an extension of HMC. It adaptively sets the step-size,  $lr$ , and the number of samples,  $N_{HMC}$ , which reduces the computation time. In Table 1, four combinations of the following hyperparameters are displayed:

- $\epsilon$  : The number of training iterations.
- $BuIn$  : The number of burn-in steps.
- $lr_0$  : The initial step-size or learning rate.
- Layers : The architecture of the network.
- $TAP$  : The Target Acceptance Probability.
- $N_{PDEs}$  : The Number of points used to evaluate the PDEs.
- Key : The random key used for the random generator.

For all the hyperparameter configurations, the maximum number of samples,  $N_{HMC}$ , is 1023. In Figure 3, the corresponding results are shown.

Table 1: Hyperparameters NUTS-B-PINN.

It	$\epsilon$	$BuIn$	$lr_0$	Layers	$TAP$	$N_{PDES}$	Key
(a)	$5e^4$	$5e^2$	$1e^{-4}$	[1, 3, 16, 20, 32, 20, 16, 3, 1]	0,80	1000	1234
(b)	$5e^4$	$5e^2$	$5e^{-4}$	[1, 3, 10, 16, 20, 32, 25, 20, 16, 10, 3, 1]	0,85	1200	12345
(c)	$5e^3$	$5e^1$	$1e^{-5}$	[1, 3, 10, 16, 20, 32, 20, 16, 10, 3, 1]	0,80	1000	20901
(d)	$5e^4$	$5e^2$	$1e^{-4}$	[1, 3, 16, 25, 16, 3, 1]	0,80	800	1234

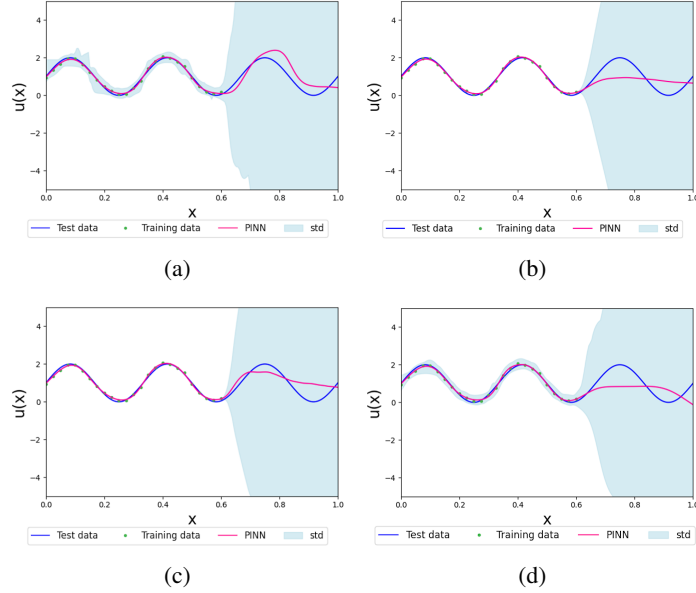


Figure 3: Results NUTS-B-PINN

### 4.3 SOLVE WITH A VI-PINN

The general architecture of VI-PINN is shown in Figure 4. The mean  $\mu_q$  is estimated with a NN. For the variance  $\sigma_q^2$ , the following the kernel is proposed

$$\sigma_q^2 = \sigma_{NN}^2 \exp[\alpha \|\mathbf{x}_i - \mathbf{x}_j\|^2]$$

with  $\sigma_{NN}^2 = \|\sigma_{NN}(\mathbf{x}_i)\sigma_{NN}(\mathbf{x}_j)\|$  where  $\sigma_{NN}(\mathbf{x})$  is a NN with input  $\mathbf{x}$ . The choice is made to take the product instead of a NN with two inputs  $(\mathbf{x}_i, \mathbf{x}_j)$ , to preserve the symmetry of the variance ( $\sigma_{NN}^2(\mathbf{x}_i, \mathbf{x}_j) = \sigma_{NN}^2(\mathbf{x}_j, \mathbf{x}_i)$ ). For this problem  $\mathcal{A}_x[\mu_q]$  and  $\mathcal{A}_x[\sigma_q^2(\mathbf{x}, \mathbf{x}')]\mathcal{A}_{x'}^T$  are given by

$$\begin{aligned} \mathcal{A}_x[\mu_q] &\equiv \nabla_x^2 \mu_q - \mu_\omega \\ \mathcal{A}_x[\sigma_q^2(\mathbf{x}, \mathbf{x}')]\mathcal{A}_{x'}^T &\equiv \nabla_x^2 \nabla_{x'}^2 \sigma_q^2(\mathbf{x}, \mathbf{x}') - \sigma_\omega^2 \delta(\mathbf{x} - \mathbf{x}') \end{aligned}$$

To reduce the training time,  $\nabla_x^2 \nabla_{x'}^2 \sigma_q^2(\mathbf{x}, \mathbf{x}')$  is calculated explicitly for this 1D problem.

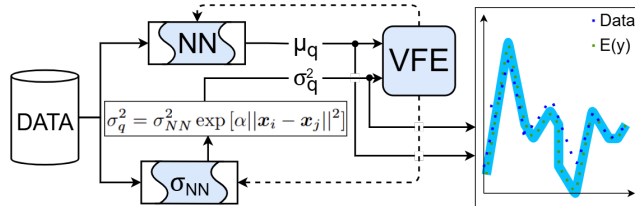


Figure 4: General architecture of a VI-PINN

The estimated mean  $\mu_q$ , and variance  $\sigma_q^2$ , are evaluated in equation (27) to find the VFE. In Table 2, the hyperparameters of the NNs for  $\mu_q$  and  $\sigma_{NN}$ , used in the VI-PINN-implementation, are shown. The intermediate results for the mean  $\mu_q$ , and the standard deviation  $\sigma_q$ , are shown respectively in Figure 5a and Figure 5b. In Figure 5c, the combined result is shown.

Table 2: Hyper parameters of the 2 NN for the VI-PINN.

NN $_{\mu_q}$	NN $_{\sigma_{NN}}$	lr	$\epsilon_{\mu_q}$	$\epsilon_{\sigma_{NN}}$	$N_{PDEs}$	$\sigma_\omega$	$\sigma_y$	$\alpha$
[1, 16, 20, 16, 1]	[1, 5, 8, 12, 12, 12, 10, 1]	$1e^{-4}$	$15e^3$	$10e^3$	700	2,0	1,5	500

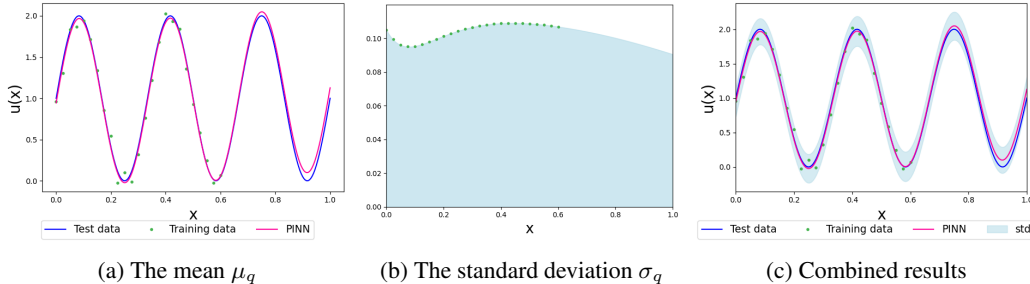


Figure 5: Results VI-PINN

#### 4.4 DISCUSS THE RESULTS

Table 3: Results from specify solving techniques

Solving technique	Training time*	MSE	AR
<i>NUTS-B-PINN</i> (a)	23 min 10,93 s	$7,93e^{-2}$	100,00%
<i>NUTS-B-PINN</i> (b)	55 min 57,72 s	$1,48e^{-1}$	41,00%
<i>NUTS-B-PINN</i> (c)	1 min 57,36 s	$9,03e^{-2}$	41,50%
<i>NUTS-B-PINN</i> (d)	8 min 49,34 s	$1,76e^{-1}$	100,00%
<i>VI-PINN</i>	4 min 39,57 s	$6,97e^{-6}$	100,00%

\*All calculations are performed on 32 CPUs (Intel(R) Core (TM) i9-14900K)

In Table 3, the training time, Means Square Error (MSE), and Acceptance Rate (AR) of each solving technique are shown.

##### 4.4.1 MEAN ESTIMATION

The two techniques both provide a good estimation for  $\mu_q$  inside the measurement range ( $x \in [0, 0.6]$ ). When looking outside the measurement range, there is a notable difference. VI-PINN can extrapolate and can yield a good solution. B-PINN can not, even when doing more iterations.

##### 4.4.2 VARIANCE ESTIMATION

For B-PINN, the variance (and standard deviation) shows a rising behaviour. Outside the measurement range ( $x \in [0.6, 1]$ ), the  $\sigma_q$  estimation exploits to high values ( $>10$ ).

When comparing to the VI-PINN, both have a similar result inside the measurement range ( $x \in [0, 0.6]$ ). When looking outside the measurement range ( $x \in [0.6, 1]$ ), a noticeable difference can be observed. The  $\sigma_q$  estimation of the VI-PINN estimation stays bounded.

This can be explained because the variance  $\sigma_q^2$  is included explicitly in the government equation of VI-PINN (see equation (28)) and in that way can be regulated. The B-PINN estimation is sample-based, hence has no explicit way of regulating it.

#### 4.4.3 COMBINED RESULT (MEAN AND STANDARD DEVIATION)

Examining the combined result of the mean  $\mu_q$  and standard deviation  $\sigma_q$  shows that the results from the VI-PINN are significantly better than those of the B-PINN, especially when extrapolating. This expresses itself in a low MSE and high AR, compared to the B-PINN estimations.

#### 4.4.4 TRAINING TIME

The training time of the B-PINN is case-dependent. When the results after an iteration have a high uncertainty, more samples are needed, and the iteration will take longer. A slight change in a hyperparameter can result in a significant increase in training time. Increasing the number of iterations does not always result in better estimations.

VI-PINN don't rely on sampling, therefore it's less sensitive to changes in hyperparameters in terms of training time, which makes it easier to tune than B-PINN, which results in a lower training time.

## 5 CONCLUSION

PINN has been shown to give good results, even with small datasets. It does so by extending the NN framework with physics, represented using PDEs or ODEs. But it has its limitations. PINNs has no uncertainty quantification built into it. In the literature, a B-PINN is proposed to solve this limitation. This is a sample-based implementation, where all parameters are assumed stochastic.

We proposed the VI-PINN as a new way of including uncertainty estimation into the PINN predictions. The VI-PINN achieves this due to its strong mathematical backbone, where VI is used to estimate the posterior in combination with VFE, which is used as a cost function in this probabilistic framework.

Comparing the B-PINN and VI-PINN, both yield an acceptable result inside the training range. Looking at the extrapolating capabilities, B-PINN falls short, a gap that VI-PINN successfully fills. The physics is firmly embedded into VI-PINN, both in the estimation of  $\mu_q$  and  $\sigma_q^2$ . This gives it the ability to give reliable results even when extrapolating. The cost function equation (27) for VI-PINN is an explicit function and therefore easier to implement and evaluate than B-PINN. This results in faster training times and results that better align with reality over the full range.

But it still has some flaws, such as a long training time and being unable to include implicit BC.

To conclude this paper, we achieve our ambition of approaching the SPDEs similarly to how PINNs approach PDEs, while the uncertainty is moved away from the parameters.

#### ACKNOWLEDGMENTS THE FUNDING

This work was realised through the financial support of the Flanders Make research project: SBO-FM-DTF-PINN.

#### CONFLICT OF INTEREST

The authors declare no potential conflict of interest.

## REFERENCES

- Saakaar Bhatnagar, Andrew Comerford, and Araz Banaeizadeh. Physics Informed Neural Networks for Modeling of 3D Flow-Thermal Problems with Sparse Domain Data. *Journal of Machine Learning for Modeling and Computing*, 5, 01 2024. doi: 10.1615/JMachLearnModelComput.2024051540.
- Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul A. Szerlip, Paul Horsfall, and Noah D. Goodman. Pyro: Deep Universal Probabilistic Programming. *J. Mach. Learn. Res.*, 20:28:1–28:6, 2019. URL <http://jmlr.org/papers/v20/18-403.html>.

- 540 William M. Bolstad and James M. Curran. Introduction to Bayesian Statistics. In *Logic, Probability,*  
541 *and Uncertainty*. New York: Wiley, 2017. 3rd ed.
- 542
- 543 Adam D. Cobb and Brian Jalaian. Scaling Hamiltonian Monte Carlo inference for Bayesian Neu-  
544 ral Networks with symmetric splitting. In Cassio de Campos and Marloes H. Maathuis (eds.),  
545 *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence*, volume  
546 161 of *Proceedings of Machine Learning Research*, pp. 675–685. PMLR, 27–30 Jul 2021. URL  
547 <https://proceedings.mlr.press/v161/cobb21a.html>.
- 548 J. Fourier. The Analytical Theory of Heat. In *The Analytical Theory of Heat*. The Univeresity, 1878.
- 549
- 550 Eric Fowler, Christopher J. McDevitt, and Subrata Roy. Physics-Informed Neural Network  
551 simulation of thermal cavity flow. *Scientific Reports*, 14(1):15203, Jul 2024. ISSN  
552 2045-2322. doi: 10.1038/s41598-024-65664-3. URL <https://doi.org/10.1038/s41598-024-65664-3>.
- 553
- 554 Ankush Ganguly and Samuel W. F. Earp. An Introduction to Variational Inference, 2021. URL  
555 <https://arxiv.org/abs/2108.13083>.
- 556
- 557 Matthew D. Hoffman and Andrew Gelman. The No-U-Turn Sampler: Adaptively Setting Path  
558 Lengths in Hamiltonian Monte Carlo, 2011. URL <https://arxiv.org/abs/1111.4246>.
- 559
- 560 Laurent Valentin Jospin, Hamid Laga, Farid Boussaid, Wray Buntine, and Mohammed Bennamoun.  
561 Hands-On Bayesian Neural Networks—A Tutorial for Deep Learning Users. *IEEE Computa-*  
562 *tional Intelligence Magazine*, 17(2):29–48, May 2022. ISSN 1556-6048. doi: 10.1109/mci.2022.  
3155327. URL <http://dx.doi.org/10.1109/MCI.2022.3155327>.
- 563
- 564 Finn Lindgren, Håvard Rue, and Johan Lindström. An explicit link between Gaussian fields and  
565 Gaussian Markov random fields: the Stochastic Partial Differential Equation approach. *Journal*  
566 *of the Royal Statistical Society Series B: Statistical Methodology*, 73(4):423–498, 2011.
- 567
- 568 R.S.H. Mah and V. Chakravarthy. Pattern recognition using artificial Neural Networks. *Com-*  
569 *puters & Chemical Engineering*, 16(4):371–377, 1992. ISSN 0098-1354. doi: [https://doi.org/10.1016/0098-1354\(92\)80054-D](https://doi.org/10.1016/0098-1354(92)80054-D). URL <https://www.sciencedirect.com/science/article/pii/009813549280054D>. Neutral network applications in chemical engineering.
- 570
- 571 Stefano Markidis. The Old and the New: Can Physics-Informed Deep-Learning Replace Tra-  
572 ditional Linear Solvers? *Frontiers in Big Data*, Volume 4 - 2021, 2021. ISSN 2624-909X.  
573 doi: 10.3389/fdata.2021.669097. URL [https://www.frontiersin.org/journals/](https://www.frontiersin.org/journals/big-data/articles/10.3389/fdata.2021.669097)  
574 [big-data/articles/10.3389/fdata.2021.669097](https://www.frontiersin.org/journals/big-data/articles/10.3389/fdata.2021.669097).
- 575
- 576 Manfred Opper. Variational Inference for Stochastic Differential Equations. *Annalen der Physik*,  
577 531(3):1800233, 2019.
- 578
- 579 K. B. Petersen and M. S. Pedersen. The Matrix Cookbook, nov 2012. URL <http://www2.compute.dtu.dk/pubdb/pubs/3274-full.html>. Version 20121115.
- 580
- 581 Du Phan, Neeraj Pradhan, and Martin Jankowiak. Composable Effects for Flexible and Accelerated  
582 Probabilistic Programming in NumPyro. *arXiv preprint arXiv:1912.11554*, 2019.
- 583
- 584 Maziar Raissi, Paris Perdikaris, and George Em Karniadakis. Physics Informed Deep Learning  
585 (part i): Data-driven Solutions of Nonlinear Partial Differential Equations, 2017. URL <https://arxiv.org/abs/1711.10561>.
- 586
- 587 Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*.  
588 MIT Press, Cambridge, MA, 2006. ISBN 9780262182539.
- 589
- 590 Simo Sarkka and Jouni Hartikainen. Infinite-dimensional Kalman Filtering approach to spatio-  
591 temporal Gaussian Process regression. In *Artificial intelligence and statistics*, pp. 993–1001.  
592 PMLR, 2012.
- 593
- 594 Dian Wu, Lei Wang, and Pan Zhang. Solving Statistical Mechanics Using Variational Autoregres-  
595 sive Networks. *Phys. Rev. Lett.*, 122:080602, Feb 2019. doi: 10.1103/PhysRevLett.122.080602.  
URL <https://link.aps.org/doi/10.1103/PhysRevLett.122.080602>.

Liu Yang, Xuhui Meng, and George Em Karniadakis. B-PINNs: Bayesian Physics-Informed Neural Networks for forward and inverse PDE problems with noisy data. *Journal of Computational Physics*, 425:109913, 2021. ISSN 0021-9991. doi: <https://doi.org/10.1016/j.jcp.2020.109913>. URL <https://www.sciencedirect.com/science/article/pii/S0021999120306872>.

## A ABBREVIATION

<b>AR</b>	Acceptance Rate
<b>BC</b>	Boundary Condition
<b>BNN</b>	Bayesian Neural Network
<b>B-PINN</b>	Bayesian PINN
<b>DTF</b>	Disturbed Thermal Field
<b>FDM</b>	Finite Difference Method
<b>FEM</b>	Finite Element Method
<b>GP</b>	Gaussian Process
<b>HMC</b>	Hamiltonian Monte Carlo
<b>NN</b>	Neural Network
<b>NUTS</b>	No-U-Turn Sampler
<b>MSE</b>	Means Square Error
<b>ODEs</b>	Ordinary Differential Equations
<b>PDEs</b>	Partial Differential Equations
<b>PINN</b>	Physics-Informed Neural Network
<b>SPDEs</b>	Stochastic Partial Differential Equations
<b>VFE</b>	Variational Free Energy
<b>VI</b>	Variational Inference
<b>VI-PINN</b>	Variational Inference Physics Informed Neural Network

## B CALCULATION OF THE TERMS OF THE VARIATIONAL FREE ENERGY

### B.1 FREE ENERGY WITHOUT MEASUREMENTS

Given:

$$\mathcal{L}_{VI}[q] = \mathbb{E}_{q(\mathbf{u})} \left[ \log \left( \frac{q(\mathbf{u})}{p(\mathbf{u})} \right) \right]$$

Calculate  $\log \left( \frac{q(\mathbf{u})}{p(\mathbf{u})} \right)$ , using

$$q(\mathbf{u}) = |2\pi\Sigma_{qq}|^{-\frac{1}{2}} \exp \left( -\frac{1}{2}(\mathbf{u} - \mu_q)^T \Sigma_{qq}^{-1} (\mathbf{u} - \mu_q) \right)$$

$$p(\boldsymbol{\omega}|\mathbf{u}) \equiv p(\mathbf{u}) = |2\pi\Sigma_{\boldsymbol{\omega}\boldsymbol{\omega}}|^{-\frac{1}{2}} \exp \left( -\frac{1}{2}(\mathbf{A}\mathbf{u} - \mu_{\boldsymbol{\omega}})^T \Sigma_{\boldsymbol{\omega}\boldsymbol{\omega}}^{-1} (\mathbf{A}\mathbf{u} - \mu_{\boldsymbol{\omega}}) \right)$$

gives

$$\log \left( \frac{q(\mathbf{u})}{p(\mathbf{u})} \right) = -\frac{1}{2} [\log(2\pi) + \log(|\Sigma_{qq}|) - \log(2\pi) - \log(|\Sigma_{\boldsymbol{\omega}\boldsymbol{\omega}}|)]$$

$$+ \frac{1}{2} [(\mathbf{A}\mathbf{u} - \mu_{\boldsymbol{\omega}})^T \Sigma_{\boldsymbol{\omega}\boldsymbol{\omega}}^{-1} (\mathbf{A}\mathbf{u} - \mu_{\boldsymbol{\omega}}) - (\mathbf{u} - \mu_q)^T \Sigma_{qq}^{-1} (\mathbf{u} - \mu_q)]$$

On the first line, the  $\log(2\pi)$  can be crossed out, and the other  $\Sigma$ 's can be taken together. Take now the expected value  $\mathbb{E}_{q(\mathbf{u})} \equiv \mathbb{E}_q$ , using Petersen & Pedersen (2012).

$$\begin{aligned} \mathcal{L}_{VI}[q] = & -\frac{1}{2} \underbrace{\mathbb{E}_q \left[ \log \left( \frac{|\Sigma_{qq}|}{|\Sigma_{\omega\omega}|} \right) \right]}_{(1)} - \frac{1}{2} \underbrace{\mathbb{E}_q \left[ (\mathbf{u} - \mu_q)^T \Sigma_{qq}^{-1} (\mathbf{u} - \mu_q) \right]}_{(2)} \\ & + \frac{1}{2} \underbrace{\mathbb{E}_q \left[ (\mathbf{A}\mathbf{u} - \mu_{\omega})^T \Sigma_{\omega\omega}^{-1} (\mathbf{A}\mathbf{u} - \mu_{\omega}) \right]}_{(3)} \end{aligned}$$

$$(1) = \log(|\Sigma_{qq}|) + \underbrace{\log(|\Sigma_{\omega\omega}|)}_{\text{IND of } q(\mathbf{u}) \rightarrow Cst.}$$

Using Equation (328) from Petersen & Pedersen (2012) to solve (2)

$$\mathbb{E}[\mathbf{c}^T \mathbf{M} \mathbf{c}] = \text{tr}(\mathbf{M} \cdot \text{cov}(\mathbf{c})) + \mu_{\mathbf{c}}^T \mathbf{M} \mu_{\mathbf{c}}$$

with

$$\mathbf{c} = (\mathbf{u} - \mu_q) \text{ and } \mathbf{M} = \Sigma_{qq}^{-1}$$

so

$$\mu_{\mathbf{c}} = 0 \text{ and } \text{cov}(\mathbf{c}) = \Sigma_{qq}$$

This results in

$$(2) = \text{tr}(\Sigma_{qq}^{-1} \Sigma_{qq}) = 1 \rightarrow Cst.$$

Use the same formula to solve (3), the result in

$$(3) = \text{tr}(\Sigma_{\omega\omega}^{-1} \mathbf{A} \Sigma_{qq} \mathbf{A}^T) + (\mathbf{A}\mu_q - \mu_{\omega})^T \Sigma_{\omega\omega}^{-1} (\mathbf{A}\mu_q - \mu_{\omega})$$

Filling in everything, equation (14) is found

$$\mathcal{L}_{VI}[q] = -\frac{1}{2} \log(|\Sigma_{qq}|) + \frac{1}{2} \text{tr}(\Sigma_{\omega\omega}^{-1} \mathbf{A} \Sigma_{qq} \mathbf{A}^T) + \frac{1}{2} (\mathbf{A}\mu_q - \mu_{\omega})^T \Sigma_{\omega\omega}^{-1} (\mathbf{A}\mu_q - \mu_{\omega}) + Cst.$$

## B.2 FREE ENERGY WITH MEASUREMENTS

Given

$$\begin{aligned} \mathcal{L}_{VI}^{\mathbf{y}}[q] &= \mathbb{E}_{q(\mathbf{u})} \left[ \log \left( \frac{q(\mathbf{u})}{p(\mathbf{u}, \mathbf{y})} \right) \right] = \mathbb{E}_{q(\mathbf{u})} \left[ \log \left( \frac{q(\mathbf{u})}{p(\mathbf{u})p(\mathbf{y}|\mathbf{u})} \right) \right] \\ &= \mathbb{E}_{q(\mathbf{u})} \left[ \log \left( \frac{q(\mathbf{u})}{p(\mathbf{u})} \right) \right] - \mathbb{E}_{q(\mathbf{u})} [\log(p(\mathbf{y}|\mathbf{u}))] \end{aligned}$$

$\mathbb{E}_{q(\mathbf{u})} \left[ \log \left( \frac{q(\mathbf{u})}{p(\mathbf{u})} \right) \right]$  is know from Appendix B.1, so only the  $\mathbb{E}_{q(\mathbf{u})} [\log(p(\mathbf{y}))]$  term must be determinate. Due to linearity of  $\mathbf{B}$ , we know that

$$p(\mathbf{y}|\mathbf{u}) = |2\pi\Sigma_{\mathbf{y}\mathbf{y}}|^{-\frac{1}{2}} \exp \left( -\frac{1}{2} (\mathbf{B}\mathbf{u} - \mu_{\mathbf{y}})^T \Sigma_{\mathbf{y}\mathbf{y}}^{-1} (\mathbf{B}\mathbf{u} - \mu_{\mathbf{y}}) \right)$$

so

$$\log(p(\mathbf{y}|\mathbf{u})) = -\frac{1}{2} (\mathbf{B}\mathbf{u} - \mu_{\mathbf{y}})^T \Sigma_{\mathbf{y}\mathbf{y}}^{-1} (\mathbf{B}\mathbf{u} - \mu_{\mathbf{y}}) + Cst.$$

Using Equation (328) from Petersen & Pedersen (2012) now for solving  $\mathbb{E}_{q(\mathbf{u})} [\log(p(\mathbf{y}|\mathbf{u}))]$

$$\mathbb{E}[\mathbf{c}^T \mathbf{M} \mathbf{c}] = \text{tr}(\mathbf{M} \cdot \text{cov}(\mathbf{c})) + \mu_{\mathbf{c}}^T \mathbf{M} \mu_{\mathbf{c}}$$

with

$$\mathbf{c} = (\mathbf{B}\mathbf{u} - \mu_{\mathbf{y}}) \text{ and } \mathbf{M} = \Sigma_{\mathbf{y}\mathbf{y}}^{-1}$$

so

$$\mu_{\mathbf{c}} = (\mathbf{B}\mu_q - \mu_{\mathbf{y}}) \text{ and } \text{cov}(\mathbf{c}) = \mathbf{B}\Sigma_{qq}\mathbf{B}^T$$

This results in

$$-\mathbb{E}_{q(\mathbf{u})} [\log(p(\mathbf{y}|\mathbf{u}))] = \frac{1}{2} \text{tr}(\Sigma_{\mathbf{y}\mathbf{y}}^{-1} \mathbf{B}\Sigma_{qq}\mathbf{B}^T) + \frac{1}{2} (\mathbf{B}\mu_q - \mu_{\mathbf{y}})^T \Sigma_{\mathbf{y}\mathbf{y}}^{-1} (\mathbf{B}\mu_q - \mu_{\mathbf{y}}) + Cst.$$