

---

# Reliable amortized variational inference with physics-based latent distribution correction

---

**Ali Siahkoohi**

School of Computational Science and Engineering,  
Georgia Institute of Technology  
alisk@gatech.edu

**Gabrio Rizzuti**

Department of Mathematics  
Utrecht University  
g.rizzuti@uu.nl

**Rafael Orozco**

School of Computational Science and Engineering,  
Georgia Institute of Technology  
rorozco@gatech.edu

**Felix J. Herrmann**

School of Computational Science and Engineering,  
Georgia Institute of Technology  
felix.herrmann@gatech.edu

## Abstract

Bayesian inference for high-dimensional inverse problems is computationally costly and requires selecting a suitable prior distribution. Amortized variational inference addresses these challenges by pretraining a neural network that approximates the posterior distribution not only for one instance of observed data, but a distribution of data pertaining to a specific inverse problem. When fed previously unseen data, the neural network—in our case a conditional normalizing flow—provides posterior samples at virtually no cost. However, the accuracy of amortized variational inference relies on the availability of high-fidelity training data, which seldom exists in geophysical inverse problems because of the Earth’s heterogeneous subsurface. In addition, the network is prone to errors if evaluated over data that is not drawn from the training data distribution. As such, we propose to increase the resilience of amortized variational inference in the presence of moderate data distribution shifts. We achieve this via a correction to the conditional normalizing flow’s latent distribution that improves the approximation to the posterior distribution for the data at hand. The correction involves relaxing the standard Gaussian assumption on the latent distribution and parameterizing it via a Gaussian distribution with an unknown mean and (diagonal) covariance. These unknowns are then estimated by minimizing the Kullback-Leibler divergence between the corrected and the (physics-based) true posterior distributions. While generic and applicable to other inverse problems, by means of a linearized seismic imaging example, we show that our correction step improves the robustness of amortized variational inference with respect to changes in the number of seismic sources, noise variance, and shifts in the prior distribution. This approach, given noisy seismic data simulated via linearized Born modeling, provides a seismic image with limited artifacts and an assessment of its uncertainty at approximately the same cost as five reverse-time migrations.

# 1 Introduction

Inverse problems involve the estimation of an unknown quantity based on noisy indirect observations. The problem is typically solved by minimizing the difference between observed and predicted data, where predicted data can be computed by modeling the underlying data generation process through a forward operator. Due to the presence of noise in the data, forward modeling errors, and the inherent nullspace of the forward operator, minimization of the data misfit alone negatively impacts the quality of the obtained solution [1]. Casting inverse problems into a probabilistic Bayesian framework allows for a more comprehensive description of their solution, where instead of finding one single solution, a distribution of solutions to the inverse problem—known as the posterior distribution—is obtained whose samples are consistent with the observed data [2]. The posterior distribution can be sampled to extract statistical information that allows for quantification of uncertainty, i.e., assessing the variability among the possible solutions to the inverse problem.

Uncertainty qualification and Bayesian inference in inverse problems often require high-dimensional posterior distribution sampling, for instance through the use of Markov chain Monte Carlo [MCMC, 3–6]. Because of their sequential nature, MCMC sampling methods require a large number of sampling steps to perform accurate Bayesian inference [7], which reduces their applicability to large-scale problems due to the high-dimensionality of the unknown and costs associated with the forward operator [4, 5, 8–14]. As an alternative, variational inference methods [15–23] approximate the posterior distribution with a surrogate and easy-to-sample distribution. By means of this approximation, sampling is turned into an optimization problem, in which the parameters of the surrogate distribution are tuned in order to minimize the divergence between the surrogate and posterior distributions. This surrogate distribution is then used for conducting Bayesian inference. While variational inference methods may have computational advantages over MCMC methods in high-dimensional inverse problems [24, 25], the resulting approximation to the posterior distribution is typically non-amortized, i.e., it is specific to the observed data used in solving the variational inference optimization problem. Thus, the variational inference optimization problem must be solved again for every new set of observations. Solving this optimization problem may require numerous iterations [19, 20], which may not be feasible in inverse problems with computationally costly forward operators, such as seismic imaging.

On the other hand, amortized variational inference [26–37] reduces Bayesian inference computational costs by incurring an up-front optimization cost for finding a surrogate conditional distribution, typically parameterized by deep neural networks [28], that approximate the posterior distribution across a family of observed data instead of being specific to a single observed dataset. This supervised learning problem involves maximization of the probability density function (PDF) of the surrogate conditional distribution over existing pairs model and data [30]. Following optimization, samples from the posterior distribution for previously unseen data may be obtained by sampling the surrogate conditional distribution, which does not require further optimization or MCMC sampling. While drastically reducing the cost of Bayesian inference, amortized variational inference can only be used for inverse problems where a dataset of model and data pairs is available that sufficiently captures the underlying joint distribution. In reality, such an assumption is rarely true in geophysical applications due to the Earth’s strong heterogeneity across geological scenarios and our lack of access to its interior [31, 38, 39]. Additionally, the accuracy of Bayesian inference with data-driven amortized variational inference methods degrades as the distribution of the data shifts with respect to pretraining data [40]. Among these shifts are changes in the distribution of noise, the number of observed data in multi-source inverse problems, and the distribution of unknowns, in other words, the prior distribution.

In this work, we leverage amortized variational inference to accelerate Bayesian inference while building resilience against data distribution shifts through an unsupervised, data-specific, a physics-based latent distribution correction method. During this process, the latent distribution of a normalizing-flow-based surrogate conditional distribution [28] is corrected to minimize the Kullback-Leibler (KL) divergence between the predicted and true posterior distributions. The invertibility of the conditional normalizing flow—a family of invertible neural networks [41]—guarantees the existence of a corrected latent distribution [42] that when “pushed forward” by the conditional normalizing flow matches the posterior distribution. During pretraining, the conditional normalizing flow learns to Gaussianize the input model and data joint samples [28], resulting in a standard Gaussian latent distribution. As a result, for slightly shifted data distributions, the conditional normalization

flow can provide samples from the posterior distribution given an “approximately Gaussian” latent distribution as input [43, 44]. Motivated by this, and to limit the costs of the latent distribution correction step, we learn a simple diagonal (elementwise) scaling and shift to the latent distribution through a physics-based objective that minimizes the KL divergence between the predicted and true posterior distributions. As with amortized variational inference, after latent distribution correction, we gain cheap access to corrected posterior samples. Besides offering computational advantages, our proposed method implicitly learns the prior distribution during conditional normalizing flow pretraining. As advocated in the literature [42, 45] learned priors have the potential to better describe the prior information when compared to generic handcrafted priors that are chosen purely for their simplicity and applicability. A schematic representation of our proposed method is shown in Figure 1.

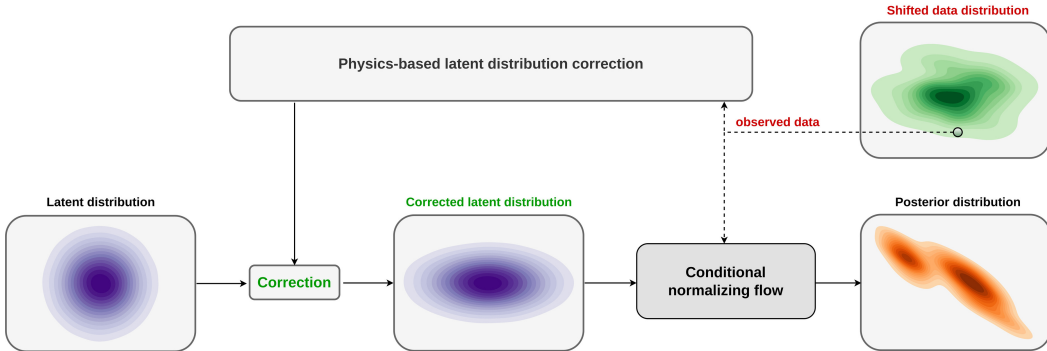


Figure 1: Schematic representation of our proposed method. We modify the standard Gaussian latent distribution of a pretrained conditional normalizing flow through a computationally cheap diagonal physics-based correction procedure to mitigate the errors due to data distribution shifts. Upon correction, the new latent samples result in corrected posterior samples when fed into the pretrained conditional normalizing flow.

## 1.1 Related work

In the context of variational inference for inverse problems, Rizzuti et al. [19], Andrieu et al. [46], Zhao et al. [47], Zhang and Curtis [48], and Zhao et al. [49] proposed a non-amortized variational inference approach to approximate posterior distributions through the use of normalizing flows. These methods do not require training data, however they require choosing a prior distribution and repeated computationally expensive evaluation of the forward operator and the adjoint of its Jacobian. Therefore, the proposed methods may prove computationally expensive when applied to inverse problems involving computationally expensive forward operators. To speed up the convergence of non-amortized variational inference, Siahkoohi et al. [31] introduces a normalizing-flow-based nonlinear preconditioning scheme. In this approach, a pretrained conditional normalizing flow capable of providing a low-fidelity approximation to the posterior distribution is used to warm-start the variational inference optimization procedure. In a related work, Kothari et al. [50] partially address challenges associated with non-amortized variational inference by learning a normalizing-flow-based prior distribution in a learned low-dimensional space via an injective network. Additionally to learning a prior, this approach also allowed non-amortized variational inference in a lower dimensional space, which could potentially have computational benefits.

Alternatively, amortized variational inference was applied by Adler and Öktem [51], Kruse et al. [28], Kovachki et al. [29], Siahkoohi and Herrmann [33], and Khorashadizadeh et al. [34] to further reduce the computational costs associated with Bayesian inference. These supervised methods learn an implicit prior distribution from training data and provide posterior samples for previously unseen data for a negligible cost due to the low cost of forward evaluation of neural networks. The success of such techniques hinges on having access to high-quality training data, including pairs of model and data that sufficiently capture the underlying model and data joint distribution. To address this limitation, Siahkoohi et al. [31] take amortized variational inference a step further by proposing a two-stage multifidelity approach where during the first stage a conditional normalizing flow is trained in the

context of amortized variational inference. To account for any potential shift in data distribution, the weights of this pretrained conditional normalizing flow are then further finetuned during an optional second stage of physics-based variational inference, which is customized for the specific imaging problem at hand. While limiting the risk of errors caused by shifts in the distribution of data, the second physics-based stage can be computationally expensive due to the high dimensionality of the weight space of conditional normalizing flows. Our work differs from the proposed method in Siahkoohi et al. [31] in that we learn to correct the latent distribution of the conditional normalizing flow, which typically has a much smaller dimensionality (approximately  $\times 90$  in our case) than the dimension of the conditional normalizing flow weight space.

The work we present is principally motivated by Asim et al. [42], which demonstrates that normalizing flows—due to their invertibility—can mitigate biases caused by shifts in the data distribution. This is achieved by reparameterizing the unknown by a pretrained normalizing flow with fixed weights while optimizing over the latent variable in order to fit the data. The reparameterization together with a Gaussian prior on the latent variable act as a regularization while the invertibility ensures the existence of a latent variable that fits the data. Asim et al. [42] exploit this property using a normalizing flow that is pretrained to capture the prior distribution associated with an inverse problem. By computing the maximum-a-posterior estimate in the latent space, Asim et al. [42], as well as Li [23] and Orozco et al. [35], limit biases originating from data distribution shifts while utilizing the prior knowledge of the normalizing flow. We extend this method by obtaining an approximation to the full posterior distribution of an inverse problem instead of a point estimate, e.g., maximum-a-posteriori.

Our work is also closely related to the non-amortized variational inference techniques presented by Whang et al. [52] and Kothari et al. [50], in which the latent distribution of a normalizing flow is altered in an unsupervised way in order to perform Bayesian inference. In contrast to our approach, these methods employ a pretrained normalizing flow that approximates the prior distribution. As a result, it is necessary to significantly alter the latent distribution in order to correct the pretrained normalizing flow to sample from the posterior distribution. In response, Whang et al. [52] and Kothari et al. [50] train a second normalizing flow aimed at learning a latent distribution that approximates the posterior distribution after passing through the pretrained normalizing flow. Our study, however, utilizes a conditional normalizing flow, which, before any corrections are applied, already approximates the posterior distribution. We argue that our approach requires a simpler correction in the latent space to mitigate biases caused by shifts in the data distribution. This is crucial when dealing with large-scale inverse problems with computationally expensive forward operators.

## 1.2 Main contributions

The main contribution of our work involves a variational inference formulation for solving probabilistic Bayesian inverse problems that leverages the benefits of data-driven learned posteriors whilst being informed by physics and data. The advantages of this formulation include

- Enhancing the solution quality of inverse problems by implicitly learning the prior distribution from the data;
- Reliably reducing the cost of uncertainty quantification and Bayesian inference; and
- Providing safeguards against data distribution shifts.

## 1.3 Outline

In the sections below, we first formulate multi-source inverse problems mathematically and cast them within a Bayesian framework. We then describe variational inference and examine how existing model and data pairs can be used to obtain an approximation to the posterior distribution that is amortized, i.e., the approximation holds over a distribution of data rather than a specific set of observations. We showcase amortized variational inference on a high-dimensional seismic imaging example in a controlled setting where we assume observed data during inference is drawn from the same distribution as training seismic data. As means to mitigate potential errors due to data distribution shifts, we introduce our proposed correction approach to amortized variational inference, which exploits the advantages of learned posteriors while reducing potential errors induced by certain data distribution shifts. Two linearized seismic imaging examples are presented, in which the distribution of the data (simulated via linearized Born modeling) is shifted by altering the forward

model and the prior distribution. These numerical experiments are intended to demonstrate the ability of the proposed latent distribution correction method to correct for errors caused by shifts in the distribution of data. Finally, we verify our proposed Bayesian inference method by conducting posterior contraction experiments.

## 2 Theory

Our purpose is to present a technique for using deep neural networks to accelerate Bayesian inference for ill-posed inverse problems while ensuring that the inference is robust with respect to data distribution shifts through the use of physics. We begin with an introduction to Bayesian inverse problems and discuss variational inference [15] as a probabilistic framework for solving Bayesian inverse problems.

### 2.1 Inverse problems

We are concerned with estimating an unknown multidimensional quantity  $\mathbf{x}^* \in \mathcal{X}$ , often referred to as the unknown model, given  $N$  noisy and indirect observed data (e.g, shot records in seismic imaging)  $\mathbf{y} = \{\mathbf{y}_i\}_{i=1}^N$  with  $\mathbf{y}_i \in \mathcal{Y}$ . Here  $\mathcal{X}$  and  $\mathcal{Y}$  denote the space of unknown models and data, respectively. The physical underlying data generation process is assumed to be encoded in forward modeling operators,  $\mathcal{F}_i : \mathcal{X} \rightarrow \mathcal{Y}$ , which relates the unknown model to the observed data via the forward model

$$\mathbf{y}_i = \mathcal{F}_i(\mathbf{x}^*) + \epsilon_i, \quad i = 1, \dots, N. \quad (1)$$

In the above expression,  $\epsilon_i$  is a vector of measurement noise, which might also include errors in the forward modeling operator. Solving ill-posed inverse problems is challenged by noise in the observed data, potential errors in the forward modeling operator, and the intrinsic nontrivial nullspace of the forward operator [1]. These challenges can lead to non-unique solutions where different estimates of the unknown model may fit the observed data equally well. Under such conditions, the use of a single model estimate ignores the intrinsic variability within inverse problem solutions, which increases the risk of overfitting the data. Therefore, not only the process of estimating  $\mathbf{x}^*$  from  $\mathbf{y}$  requires regularization, but it also calls for a statistical inference framework that allows us to characterize the variability among the solutions by quantifying the solution uncertainty [2].

### 2.2 Bayesian inference for solving inverse problems

To systematically quantify the uncertainty, we cast the inverse problem into a Bayesian framework [2]. In this framework, instead of having a single estimate of the unknown, the solution is characterized by a probability distribution over the solution space  $\mathcal{X}$  that is conditioned on data, namely the posterior distribution. This conditional distribution, denoted by  $p_{\text{post}}(\mathbf{x} | \mathbf{y})$ , can according to the Bayes' rule be written as follows:

$$p_{\text{post}}(\mathbf{x} | \mathbf{y}) = \frac{p_{\text{like}}(\mathbf{y} | \mathbf{x}) p_{\text{prior}}(\mathbf{x})}{p_{\text{data}}(\mathbf{y})}. \quad (2)$$

which equivalently can be expressed as

$$\begin{aligned} -\log p_{\text{post}}(\mathbf{x} | \mathbf{y}) &= -\sum_{i=1}^N \log p_{\text{like}}(\mathbf{y}_i | \mathbf{x}) - \log p_{\text{prior}}(\mathbf{x}) + \log p_{\text{data}}(\mathbf{y}) \\ &= \frac{1}{2\sigma^2} \sum_{i=1}^N \|\mathbf{y}_i - \mathcal{F}_i(\mathbf{x})\|_2^2 - \log p_{\text{prior}}(\mathbf{x}) + \text{const}, \end{aligned} \quad (3)$$

in case the observed data ( $\mathbf{y}_i$ ) are independent conditioned on the unknown model  $\mathbf{x}$ . In equations 2 and 3, the likelihood function  $p_{\text{like}}(\mathbf{y} | \mathbf{x})$  quantifies how well the predicted data fits the observed data given the PDF of the noise distribution. For simplicity, we assume the distribution of the noise is a zero-mean Gaussian distribution with covariance  $\sigma^2 \mathbf{I}$  but other choices can be incorporated. The prior distribution  $p_{\text{prior}}(\mathbf{x})$  encodes prior beliefs on the unknown quantity, which can also be interpreted as a regularizer for the inverse problem. Finally,  $p_{\text{data}}(\mathbf{y})$  denotes the data PDF, which is a normalization constant that is independent of  $\mathbf{x}$ .

Acquiring statistical information regarding the posterior distribution requires access to samples from the posterior distribution. Sampling the posterior distribution, commonly achieved via MCMC

[3] or variational inference techniques [15], is computationally costly in high-dimensional inverse problems due to the costs associated with many needed evaluations of the forward operator [2, 4, 5, 11, 12, 24, 25, 53–56]. For multi-source inverse problem the costs are especially high as evaluating the likelihood function involves  $N$  forward operator evaluations (equation 3). Stochastic gradient Langevin dynamics [SGLD; 6, 9, 57] alleviates the need to evaluate the likelihood for all the  $N$  forward operators by allowing for stochastic approximations to the likelihood, i.e., evaluating the likelihood over randomly selected indices  $i \in \{1, \dots, N\}$ . While SGLD can provably provide accurate posterior samples with more favorable computational costs [9], due to the sequential nature of MCMC methods, SGLD still requires numerous iterations to fully traverse the probability space [7], which is computationally challenging in large-scale multi-source inverse problems. In the next section, we introduce variational inference as an alternative Bayesian inference method that has the potential to scale better than MCMC-methods in inverse problems with costly forward operators [24, 25].

### 2.3 Variational inference

As an alternative to MCMC-based methods, variational inference methods [15] reduce the problem of sampling from the posterior distribution  $p_{\text{post}}(\mathbf{x} | \mathbf{y})$  to an optimization problem. The optimization problem involves approximating the posterior PDF via the PDF of a tractable surrogate distribution  $p_\phi(\mathbf{x})$  with parameters  $\phi$  by minimizing a divergence (read “distance”) between  $p_\phi(\mathbf{x})$  and  $p_{\text{post}}(\mathbf{x} | \mathbf{y})$  with respect to surrogate distribution parameters  $\phi$ . This optimization problem can be solved approximately, which allows for trading off computational cost for accuracy [15]. After optimization, we gain access to samples from the posterior distribution by sampling  $p_\phi(\mathbf{x})$  instead, which does not involve forward operator evaluations.

Due to its simplicity and connections to the maximum likelihood principle [58], we formulate variational inference via the Kullback-Leibler (KL) divergence. The KL divergence can be explained as the cross-entropy of  $p_{\text{post}}(\mathbf{x} | \mathbf{y})$  relative to  $p_\phi(\mathbf{x})$  minus the entropy of  $p_\phi(\mathbf{x})$ . This definition describes the reverse KL divergence, denoted by  $\mathbb{KL}(p_\phi(\mathbf{x}) || p_{\text{post}}(\mathbf{x} | \mathbf{y}))$ , which is not equal to the forward KL divergence,  $\mathbb{KL}(p_{\text{post}}(\mathbf{x} | \mathbf{y}) || p_\phi(\mathbf{x}))$ . This non-symmetry in KL divergence leads to different computational and approximation properties during variational inference, which we describe in detail in the following sections. We will first describe the reverse KL divergence, followed by the forward KL divergence. Finally, we will describe normalizing flows as a way of parameterized surrogate distributions to facilitate variational inference.

#### 2.3.1 Non-amortized variational inference

The reverse KL divergence is the common choice for formulating variational inference [19, 46–49] in which the physically-informed posterior density guides the optimization over  $\phi$ . The reverse KL divergence can be mathematically stated as

$$\mathbb{KL}(p_\phi(\mathbf{x}) || p_{\text{post}}(\mathbf{x} | \mathbf{y}_{\text{obs}})) = \mathbb{E}_{\mathbf{x} \sim p_\phi(\mathbf{x})} [-\log p_{\text{post}}(\mathbf{x} | \mathbf{y}_{\text{obs}}) + \log p_\phi(\mathbf{x})], \quad (4)$$

where  $\mathbf{y}_{\text{obs}} \sim p_{\text{data}}(\mathbf{y})$  refers to a specific single observed data.  $\mathbf{x}$  in the right hand side of the expression in equation 4 is a random variable obtained by sampling the surrogate distribution  $p_\phi(\mathbf{x})$ , over which we evaluate the expectation. Variational inference using the reverse KL divergence involves minimizing equation 4 with respect to  $\phi$  during which the logarithm of the posterior PDF is approximated by the logarithm of the surrogate PDF, when evaluated over samples from the surrogate distribution. By expanding the negative-log posterior density via Bayes’ rule (equation 3), we write the non-amortized variational inference optimization problem as

$$\phi^* = \arg \min_{\phi} \mathbb{E}_{\mathbf{x} \sim p_\phi(\mathbf{x})} \left[ \frac{1}{2\sigma^2} \sum_{i=1}^N \|\mathbf{y}_{\text{obs},i} - \mathcal{F}_i(\mathbf{x})\|_2^2 - \log p_{\text{prior}}(\mathbf{x}) + \log p_\phi(\mathbf{x}) \right]. \quad (5)$$

The expectation in the above equation is approximated with a sample mean over samples drawn from  $p_\phi(\mathbf{x})$ . The optimization problem in equation 5 can be solved using stochastic gradient descent and its variants [59–62] where at each iteration the objective function is evaluated over a batch of samples drawn from  $p_\phi(\mathbf{x})$  and randomly selected (without replacement) indices  $i \in \{1, \dots, N\}$ . To solve this optimization problem, there are two considerations to take into account. First consideration involves the tractable computation of the surrogate PDF and its gradient with respect to  $\phi$ . As described in the

following sections, normalizing flows [17], which are a family of specially designed invertible neural networks [41], facilitate the computation of these quantities via the change-of-variable formula in probability distributions [63]. The second consideration involves differentiating (with respect to  $\phi$ ) the expectation (sample mean) operation in equation 5. Evaluating this expectation requires sampling from the surrogate distribution  $p_\phi(\mathbf{x})$ , which depends  $\phi$ . Differentiating through the sampling procedure from the surrogate distribution  $p_\phi(\mathbf{x})$  can be facilitated through the reparameterization trick [64]. In this approach sampling from  $p_\phi(\mathbf{x})$  is interpreted as passing latent samples  $\mathbf{z} \in \mathcal{Z}$  from a simple base distribution, such as standard Gaussian distribution, through a parametric function parameterized by  $\phi$  [64]. With this interpretation, the expectation over  $p_\phi(\mathbf{x})$  can be computed over the latent distribution instead, which does not depend on  $\phi$ , followed by a mapping of latent samples through the parametric function. This process enables computing the gradient of the expression in equation 5 with respect to  $\phi$  [64].

Following optimization,  $p_{\phi^*}(\mathbf{x})$  provides unlimited samples from the posterior distribution—virtually for free. While there are indications that this approach can be computationally favorable compared to MCMC sampling methods [24, 25], each iteration during optimization problem 5 involves evaluating the forward operator and the adjoint of its Jacobian, which can be computationally costly depending on  $N$  and the number of iterations required to solve 5. In addition, and more importantly, this approach is non-amortized—i.e., the resulting surrogate distribution  $p_{\phi^*}(\mathbf{x})$  approximates the posterior distribution for the specific data  $\mathbf{y}_{\text{obs}}$  that is used to solve optimization problem 5. This necessitates the optimization problem to be solved again for a new instance of the inverse problem with different data. In the next section, we introduce an amortized variational inference approach that addresses these limitations.

### 2.3.2 Amortized variational inference

Similarly to reverse KL divergence, forward KL divergence involves calculating the difference between the logarithms of the surrogate PDF and the posterior PDF. In contrast to reverse KL divergence, however, to compute the forward KL divergence the PDFs are evaluated over samples from the posterior distribution rather than the surrogate distribution samples (see equation 4). The forward KL divergence can be written as follows

$$\mathbb{KL}(p_{\text{post}}(\mathbf{x} | \mathbf{y}) || p_\phi(\mathbf{x})) = \mathbb{E}_{\mathbf{x} \sim p_{\text{post}}(\mathbf{x} | \mathbf{y})} \left[ -\log p_\phi(\mathbf{x}) + \log p_{\text{post}}(\mathbf{x} | \mathbf{y}) \right]. \quad (6)$$

Following the expression above, it is infeasible to evaluate the forward KL divergence in inverse problems as it requires access to samples from the posterior distribution—the samples that we are ultimately after and do not have access to. However, the average (over data) forward KL divergence can be computed using available model and data pairs in the form of samples from the joint distribution  $p(\mathbf{x}, \mathbf{y})$ . This involves integrating (marginalizing) the forward KL divergence over existing data  $\mathbf{y} \sim p_{\text{data}}(\mathbf{y})$ :

$$\begin{aligned} & \mathbb{E}_{\mathbf{y} \sim p_{\text{data}}(\mathbf{y})} \left[ \mathbb{KL}(p_{\text{post}}(\mathbf{x} | \mathbf{y}) || p_\phi(\mathbf{x})) \right] \\ &= \mathbb{E}_{\mathbf{y} \sim p_{\text{data}}(\mathbf{y})} \mathbb{E}_{\mathbf{x} \sim p_{\text{post}}(\mathbf{x} | \mathbf{y})} \left[ -\log p_\phi(\mathbf{x} | \mathbf{y}) + \underbrace{\log p_{\text{post}}(\mathbf{x} | \mathbf{y})}_{\text{constant w.r.t. } \phi} \right] \\ &= \iint \underbrace{p_{\text{data}}(\mathbf{y}) p_{\text{post}}(\mathbf{x} | \mathbf{y})}_{=p(\mathbf{x}, \mathbf{y})} \left[ -\log p_\phi(\mathbf{x} | \mathbf{y}) \right] d\mathbf{x} d\mathbf{y} + \text{const} \\ &= \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p(\mathbf{x}, \mathbf{y})} \left[ -\log p_\phi(\mathbf{x} | \mathbf{y}) \right] + \text{const}. \end{aligned} \quad (7)$$

In the above expression  $p_\phi(\mathbf{x} | \mathbf{y})$  represents a surrogate conditional distribution that approximates the posterior distribution for any data  $\mathbf{y} \sim p_{\text{data}}(\mathbf{y})$ . The third line in equation 7 is the result of applying the chain rule of PDFs<sup>1</sup>. By minimizing the average KL divergence we obtain the following amortized variational inference objective:

$$\begin{aligned} \phi^* &= \arg \min_{\phi} \mathbb{E}_{\mathbf{y} \sim p_{\text{data}}(\mathbf{y})} \left[ \mathbb{KL}(p_{\text{post}}(\mathbf{x} | \mathbf{y}) || p_\phi(\mathbf{x} | \mathbf{y})) \right] \\ &= \arg \min_{\phi} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p(\mathbf{x}, \mathbf{y})} \left[ -\log p_\phi(\mathbf{x} | \mathbf{y}) \right]. \end{aligned} \quad (8)$$

---

<sup>1</sup> $p(x, y) = p(x | y) p(y), \forall x \in \mathcal{X}, y \in \mathcal{Y}$ .

The above optimization problem represent a supervised learning framework for obtaining fully-learned posteriors using existing pairs of model and data. The expectation is approximated with a sample mean over available model and data joint samples. Note that this method does not impose any explicit assumption on the noise distribution (see equation 3), and the information about the forward model is implicitly encoded in the model and data pairs. As a result, this formulation is an instance of likelihood-free simulation-based inference methods [65, 66] that allows us to approximate the posterior distribution for previously unseen data as,

$$p_{\phi^*}(\mathbf{x} \mid \mathbf{y}_{\text{obs}}) \approx p_{\text{post}}(\mathbf{x} \mid \mathbf{y}_{\text{obs}}), \quad \forall \mathbf{y}_{\text{obs}} \sim p_{\text{data}}(\mathbf{y}). \quad (9)$$

Equation 9 holds for previously unseen data drawn from  $p_{\text{data}}(\mathbf{y})$  provided that the optimization problem 8 is solved accurately [28, 40], i.e.,  $\mathbb{E}_{\mathbf{y} \sim p_{\text{data}}(\mathbf{y})} [\mathbb{KL}(p_{\text{post}}(\mathbf{x} \mid \mathbf{y}) \parallel p_{\phi^*}(\mathbf{x} \mid \mathbf{y}))] = 0$ . Following an one-time upfront cost of training, equation 9 can be used to sample the posterior distribution with no additional forward operator evaluations. While computationally cheap, the accuracy of the amortized variational inference approach in equation 8 is directly linked to the quality and quantity of model and data pairs used during optimization [65]. This raises questions regarding the reliability of this approach in domains that sufficiently capturing the underlying joint model and data distribution is challenging, e.g., in geophysical applications due to the Earth’s strong heterogeneity across geological scenarios and our lack of access to its interior [31, 38, 39]. To increase the resilience of amortized variational inference when faced with data distribution shifts, e.g., changes in the forward model or prior distribution, we propose a latent distribution correction to physically inform the inference. Before describing our proposed physics-based latent distribution correction approach, we introduce conditional normalizing flows [28] to parameterize the surrogate conditional distribution for amortized variational inference.

## 2.4 Conditional normalizing flows for amortized variational inference

To limit the computational cost of amortized variational inference, both during optimization and inference, it is imperative that the surrogate conditional distribution be able to: (1) approximate complex distributions, i.e., it should have a high representation power, which is required to represent possibly multi-modal distributions; (2) support cheap density estimation, which involves computing the density  $p_{\phi}(\mathbf{x} \mid \mathbf{y})$  for given  $\mathbf{x}$  and  $\mathbf{y}$ ; and (3) permit fast sampling from  $p_{\phi}(\mathbf{x} \mid \mathbf{y})$  for cheap posterior sampling during inference. These characteristics are provided by conditional normalizing flows [28], which are a family of invertible neural networks [41] that are capable of approximating complex conditional distributions [67, 68].

A conditional normalizing flows—in the context of amortized variational inference—aims to map input samples  $\mathbf{z}$  from a latent standard multivariate Gaussian distribution  $\mathcal{N}(\mathbf{z} \mid \mathbf{0}, \mathbf{I})$  to samples from the posterior distribution given the observed data  $\mathbf{y} \sim p_{\text{data}}(\mathbf{y})$  as an additional input. This nonlinear mapping can formally be stated as  $f_{\phi}^{-1}(\cdot; \mathbf{y}) : \mathcal{Z} \rightarrow \mathcal{X}$ , with  $f_{\phi}^{-1}(\mathbf{z}; \mathbf{y})$  being the inverse of the conditional normalizing flow with respect to its first argument. Due to the low computational cost of evaluating invertible neural networks in reverse [41], using conditional normalizing flows as a surrogate conditional distribution  $p_{\phi}(\mathbf{x} \mid \mathbf{y})$  allows for extremely fast sampling from  $p_{\phi}(\mathbf{x} \mid \mathbf{y})$ . In addition to low-cost sampling, the invertibility of conditional normalizing flows permits straightforward and cheap estimation of the density  $p_{\phi}(\mathbf{x} \mid \mathbf{y})$ . This allows for tractable amortized variational inference via equation 8 through the following change-of-variable formula in probability distributions [63],

$$p_{\phi}(\mathbf{x} \mid \mathbf{y}) = \mathcal{N}(f_{\phi}(\mathbf{x}; \mathbf{y}) \mid \mathbf{0}, \mathbf{I}) \left| \det \nabla_{\mathbf{x}} f_{\phi}(\mathbf{x}; \mathbf{y}) \right|, \quad \forall \mathbf{x}, \mathbf{y} \sim p(\mathbf{x}, \mathbf{y}). \quad (10)$$

In the above formula,  $\mathcal{N}(f_{\phi}(\mathbf{x}; \mathbf{y}) \mid \mathbf{0}, \mathbf{I})$  represents the PDF for a multivariate standard Gaussian distribution evaluated at  $f_{\phi}(\mathbf{x}; \mathbf{y})$ . Thanks to the special design of invertible neural networks [41], density estimation via equation 10 is cheap since evaluating the conditional normalizing flow and the determinant of its Jacobian  $\det \nabla_{\mathbf{x}} f_{\phi}(\mathbf{x}; \mathbf{y})$  are almost free of cost. Given the expression for  $p_{\phi}(\mathbf{x} \mid \mathbf{y})$  in equation 10, we derive the following training objective for amortized conditional normalizing flows:

$$\begin{aligned} \phi^* &= \arg \min_{\phi} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p(\mathbf{x}, \mathbf{y})} [-\log p_{\phi}(\mathbf{x} \mid \mathbf{y})] \\ &= \arg \min_{\phi} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p(\mathbf{x}, \mathbf{y})} \left[ \frac{1}{2} \|f_{\phi}(\mathbf{x}; \mathbf{y})\|_2^2 - \log \left| \det \nabla_{\mathbf{x}} f_{\phi}(\mathbf{x}; \mathbf{y}) \right| \right]. \end{aligned} \quad (11)$$



In the above objective, the  $\ell_2$ -norm follows from a standard Gaussian distribution assumption on the latent variable, i.e., the output of the normalizing flow. The second term quantifies the relative change of density volume [papamakarios2021] and can be interpreted as an entropy regularization of  $p_\phi(\mathbf{x} | \mathbf{y})$ , which prevents the conditional normalizing flow from converging to solutions, e.g.,  $f_\phi(\mathbf{x}; \mathbf{y}) := \mathbf{0}$ . Due to the particular design of invertible networks [28, 41], computing the gradient of  $\det \nabla_{\mathbf{x}} f_\phi(\mathbf{x}; \mathbf{y})$  has a negligible extra cost. Figure 2 illustrates the pretraining phase as a schematic.

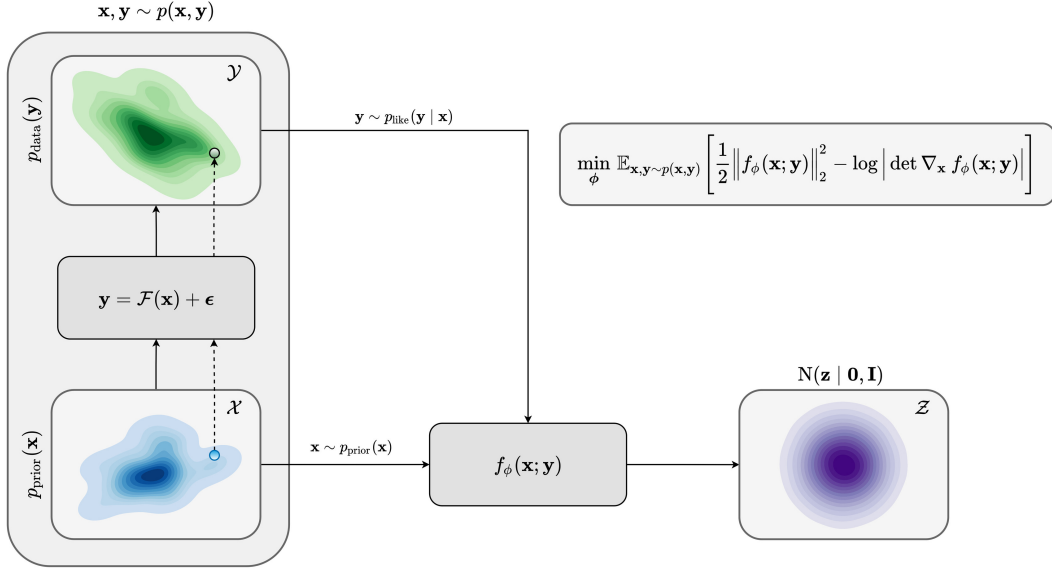


Figure 2: A schematic representation of pretraining conditional normalizing flows in the context of amortized variational inference. During pretraining, joint model and data joint samples  $\mathbf{x}, \mathbf{y} \sim p(\mathbf{x}, \mathbf{y})$  from the training dataset and are fed to the conditional normalizing flow. The training objective (equation 11) enforces the conditional normalizing flow to Gaussianize its input.

After training, given a previously unseen observed data  $\mathbf{y}_{\text{obs}} \sim p_{\text{data}}(\mathbf{y})$  we sample from the posterior distribution using the inverse of the conditional normalizing flow. We achieve this by feeding latent samples  $\mathbf{z} \sim \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})$  to the conditional normalizing flow’s inverse  $f_\phi^{-1}(\mathbf{z}; \mathbf{y}_{\text{obs}})$  while conditioning on the observed data  $\mathbf{y}_{\text{obs}}$ ,

$$f_\phi^{-1}(\mathbf{z}; \mathbf{y}_{\text{obs}}) \sim p_{\text{post}}(\mathbf{x} | \mathbf{y}_{\text{obs}}), \quad \mathbf{z} \sim \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I}). \quad (12)$$

This step is illustrated in Figure 3. As the process above does not involve forward operator evaluations, sampling with pretrained conditional normalizing flows is fast once an upfront cost of amortized variational inference is incurred. In the next section, we apply the above amortized variational inference to a seismic imaging example in a controlled setting in which we assume no data distribution shifts during inference.

### 3 Validating amortized variational inference

The objective of this example is to apply amortized variational inference to the high-dimensional seismic imaging problem. We show that a relatively good pretrained conditional normalizing flow within the context of amortized variational inference can be used to provide approximate posterior samples for previously unseen seismic data that is drawn from the same distribution as training seismic data. We begin by introducing seismic imaging and describe challenges with Bayesian inference in this problem.

#### 3.1 Seismic imaging

We are concerned with constructing an image of the Earth’s subsurface using indirect surface measurements that record the Earth’s response to synthetic sources being fired on the surface. The nonlinear relationship between these measurements, known as shot records, and the squared-slowness

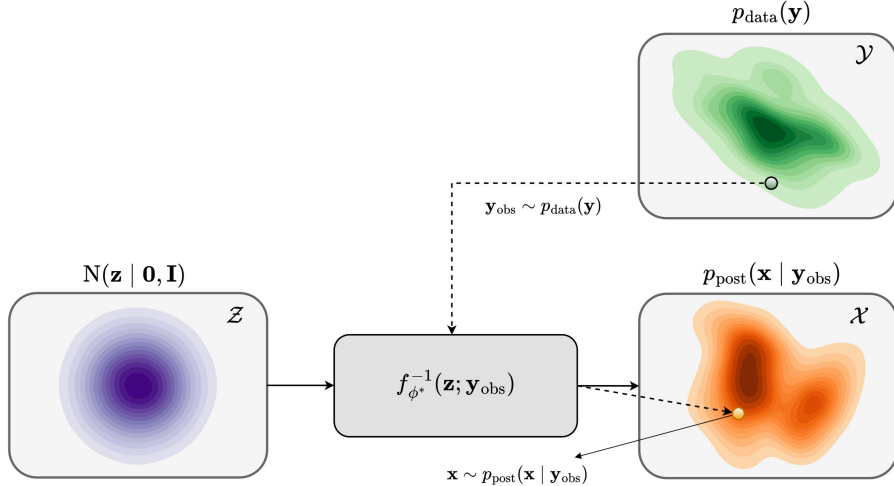


Figure 3: A schematic representation of posterior sampling with pretrained conditional normalizing flows. To sample from  $p_{\text{post}}(\mathbf{x} | \mathbf{y}_{\text{obs}})$ , the observed data and latent samples are fed to the conditional normalizing flow’s inverse  $f_{\phi}^{-1}(\mathbf{z}; \mathbf{y}_{\text{obs}})$ . Each latent sample realization results in a realization of the posterior distribution.

model of the Earth’s subsurface is governed by the wave equation. By linearizing this nonlinear relation, seismic imaging aims to estimate the short-wavelength component of the Earth’s subsurface squared-slowness model. In its simplest acoustic form, the linearization with respect to the slowness model—around a known, smooth background squared slowness model  $\mathbf{m}_0$ —leads to the following linear forward problem:

$$\mathbf{d}_i = \mathbf{J}(\mathbf{m}_0, \mathbf{q}_i) \delta \mathbf{m}^* + \epsilon_i, \quad i = 1, \dots, N. \quad (13)$$

We invert the above forward model to estimate the ground truth seismic image  $\delta \mathbf{m}^*$  from  $N$  processed (linearized) shot records  $\{\mathbf{d}_i\}_{i=1}^N$  where  $\mathbf{J}(\mathbf{m}_0, \mathbf{q}_i)$  represents the linearized Born scattering operator [69]. This operator is parameterized by the source signature  $\mathbf{q}_i$  and the smooth background squared-slowness model  $\mathbf{m}_0$ . Noise is denoted by  $\epsilon_i$ , and represents measurement noise and linearization errors. While amortized variational inference does not require knowing the closed form expression of the noise density to simulate pairs of data and model (e.g., it is sufficient to be able to simulate noise instances), for simplicity we assume the noise distribution is a zero-centered Gaussian distribution with known covariance  $\sigma^2 \mathbf{I}$ . Due to the presence of shadow zones and noisy finite-aperture shot data, wave-equation based linearized seismic imaging (in short seismic imaging for the purposes of this paper) corresponds to solving an inconsistent and ill-conditioned linear inverse problem [70–72]. To avoid the risk of overfitting the data and to quantify uncertainty, we cast the seismic imaging problem into a Bayesian inverse problem [2].

To address the challenge of Bayesian inference in this high-dimensional inverse problem, we adhere to our amortized variational inference framework. Within this approach, for an one-time upfront cost of training a conditional normalizing flow, we get access to posterior samples for previously unseen observed data that are drawn from the same distribution as the distribution of training seismic data. This includes data acquired in areas of the Earth with similar geologies, e.g., in neighboring surveys. In addition, in our framework no explicit prior density function needs to be chosen as the conditional normalizing flow learns the prior distribution during pretraining from the collection of seismic images in the training dataset. The implicitly learned prior distribution by the conditional normalizing flow minimizes the risk of negatively biasing the outcome of Bayesian inference by using overly simplistic priors. In the next section, we describe the setup for our amortized variational inference for seismic imaging.

### 3.1.1 Acquisition geometry

To mimic the complexity of real seismic images, we propose a “quasi”-real data example in which we generate synthetic data by applying the linearized Born scattering operator to 4750 2D sections with

size  $3075 \text{ m} \times 5120 \text{ m}$  extracted from the shallow section of the Kirchhoff migrated [Parihaka-3D](#) field dataset [73, 74]. We consider a 12.5 m vertical and 20 m horizontal grid spacing, and we augment an artificial 125 m water column on top of these images. We parameterize the linearized Born scattering operator via a fictitious background squared-slowness model, derived from the Kirchhoff migrated images. To ensure good coverage, we simulate 102 shot records with a source spacing of 50 m. Each shot is recorded for two seconds with 204 fixed receivers sampled at 25 m spread on top of the model. The source is a Ricker wavelet with a central frequency of 30 Hz. To mimic a more realistic imaging scenario, we add band-limited noise to the shot records, where the noise is obtained by filtering white noise with the source wavelet (Figure 5b).

### 3.1.2 Training configuration

Casting seismic imaging into amortized variational inference, as described in this paper, is hampered by the high-dimensionality of the data due to the multi-source nature of this inverse problem. To avoid computational complexities associated with directly using  $N$  shot records as input to the conditional normalizing flow, we choose to condition the conditional normalizing flow on the reverse-time migrated image, which can be estimated by applying the adjoint of the linearized Born scattering operator to the shot records,

$$\delta \mathbf{m}_{\text{RTM}} = \sum_{i=1}^N \mathbf{J}(\mathbf{m}_0, \mathbf{q}_i)^\top \mathbf{d}_i. \quad (14)$$

While  $\mathbf{d}_i$  in the above expression is defined according to the linearized forward model in equation 13, which does not involve linearization errors, our method can handle observed data simulated from wave-equation based nonlinear forward modeling. Conditioning on the reverse-time migrated image and not on the shot records directly may result in learning an approximation to the true posterior distribution [75]. While technique from statistics involving learned summary functions [30, 40] can reduce the dimensionality of the observed data, we propose to limit the Bayesian inference bias induced by conditioning on the reverse-time migrated image via our physics-based latent variable correction approach. We leave utilizing summary functions in the context of seismic imaging Bayesian inference to future work.

To create training pairs,  $(\delta \mathbf{m}^{(i)}, \delta \mathbf{m}_{\text{RTM}}^{(i)})$ ,  $i = 1, \dots, 4750$ , we first simulate (see Figure 2) noisy seismic data according to the above-mentioned acquisition design for all extracted seismic images  $\delta \mathbf{m}^{(i)}$  from shallow sections of the imaged Parihaka dataset. Next, we compute  $\delta \mathbf{m}_{\text{RTM}}^{(i)}$  by applying reverse-time-migration to the observed data for each image  $\delta \mathbf{m}^{(i)}$ . As for the conditional normalizing flow architecture, we follow Kruse et al. [28] and use hierarchical normalizing flows due to their increased expressiveness when compared to conventional invertible architectures [41]. The expressive power of hierarchical normalizing flows is a result of applying a series of conventional invertible layers [41] to different scales of the input in a hierarchical manner (refer to Kruse et al. [28] for a schematic representation of the architecture). This leads to an invertible architecture with a dense Jacobian [28] that is capable of representing complicated bijective transformations. We train this conditional normalizing flow on the pairs  $(\delta \mathbf{m}^{(i)}, \delta \mathbf{m}_{\text{RTM}}^{(i)})$ ,  $i = 1, \dots, 4750$  according to the objective function in equation 11 with the Adam stochastic optimization method [62] with a batchsize of 16 for one thousand passes over the training dataset (epochs). We use an initial stepsize of  $10^{-4}$  and decrease it after each epoch until reaching the final stepsize of  $10^{-6}$ . To monitor overfitting, we evaluate the objective function at the end of every epoch over random subsets of the validation set, consisting of 530 seismic images extracted from the shallow sections of the imaged Parihaka dataset and the associated reverse-time migrated images. As illustrated in Figure 4, the training and validation objective values exhibit a decreasing trend, which suggests no overfitting. We stopped the training after one thousand epochs due to a slowdown in the decrease of the training and validation objective values.

## 3.2 Results and observations

Following training, the pretrained conditional normalizing flow is able to produce samples from the posterior distribution for seismic data not used in training. These samples resemble different regularized (via the learned prior) least-squares migration images that explain the observed data. To demonstrate this, we simulate seismic data for a previously unseen perturbation model using the

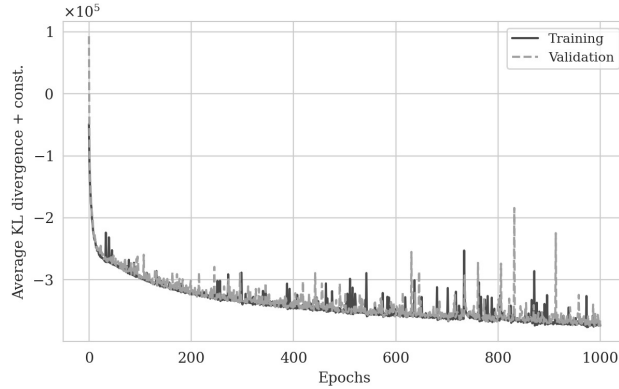


Figure 4: Training and validation objective values as a function of epochs. The validation objective value is computed over randomly selected batched of the validation set at the end of each epoch.

forward model 13 with the same noise variance. Figure 5 shows an example of a single noise-free (Figure 5a) and noisy (Figure 5b) shot record for one of 102 sources.

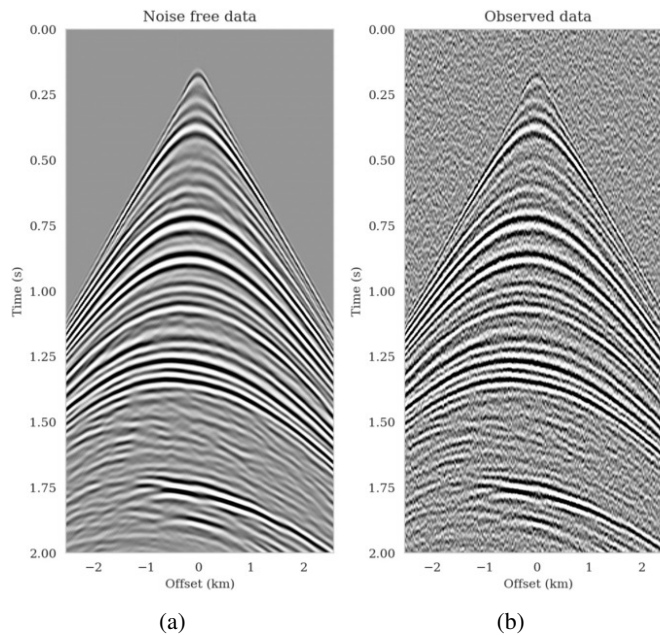
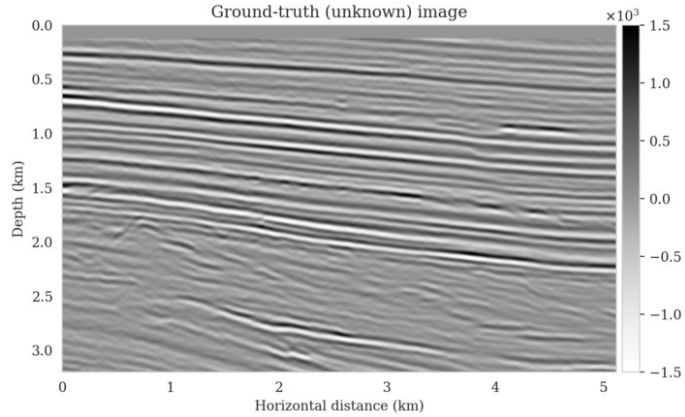


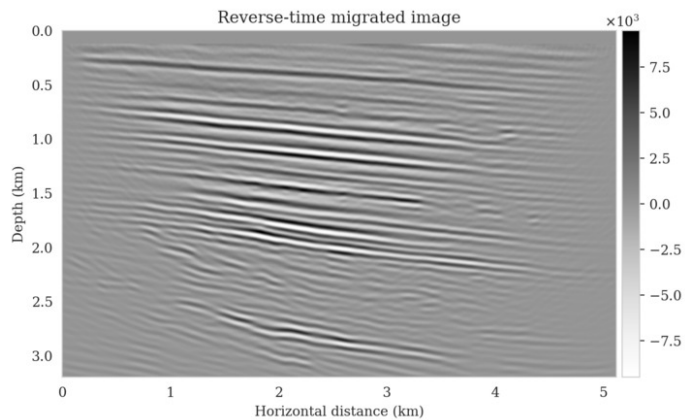
Figure 5: A shot record generated from an image extracted from the Parihaka dataset. (a) Noise-free linearized data. (b) Linearized data with bandwidth-limited noise.

We perform reverse-time migration to obtain the necessary input for the conditional normalizing flow to obtain posterior samples. We show the ground-truth seismic image (to be estimated) and the resulting reverse-time migrated image in Figures 6a and 6b, respectively. Clearly, the reverse-time migrated image has grossly wrong amplitudes, and more importantly, due to limited-aperture shot data, the edges of the image are not well illuminated.

We obtain one thousand posterior samples by providing the reverse-time migrated image and latent samples drawn from the standard Gaussian distribution to the pretrained conditional normalizing flow (equation 9). This process is fast as it does not require any forward operator evaluations. To illustrate the variability among the posterior samples, we show six of them in Figure 7. As shown in Figure 7, these image samples have amplitudes in the same range as the ground-truth image and better predict the reflectors at the edges of the image compared to the reverse-time migration image (Figure 6b). In



(a)



(b)

Figure 6: Amortized variational inference testing phase setup. (a) High-fidelity ground-truth image. (b) Reverse-time migrated image with SNR  $-12.17$  dB.

addition, the posterior samples indicate improved imaging in deep regions, which is typically more difficult due to the placement of the sources and receiver near the surface.

Samples from the posterior provide access to useful statistical information including approximations to moments of the distribution such as the mean and pointwise standard deviation (Figure 8). We compute the mean of the posterior samples to obtain the conditional mean estimate, i.e., the expected value of the posterior distribution. This estimate is depicted in Figure 8a. From Figure 8a, we observe that the overall amplitudes are well recovered by the conditional mean estimate, which includes partially recovered reflectors in badly illuminated areas close to the boundaries. Although the reconstructions are not perfect, they significantly improve upon the reverse-time migrations estimate. We did not observe a significant increase in the signal-to-noise ratio (SNR) of the conditional mean estimate when more than one thousand samples from the posterior are drawn. We use the one thousand samples to also estimate the pointwise standard deviation (Figure 8b), which serves as an assessment of the uncertainty. To avoid bias from strong amplitudes in the estimated image, we also plot the stabilized division of the standard deviation by the envelope of the conditional mean in Figure 8c. As expected, the pointwise standard deviation in Figures 8b and 8c indicate that we have the most uncertainty in areas of complex geology—e.g., near channels and tortuous reflectors, and in areas with a relatively poor illumination (deep and close to boundaries). The areas with large uncertainty align well with difficult-to-image parts of the model. The normalized pointwise standard deviation (Figure 8c) aims to visualize an amplitude-independent assessment of uncertainty, which indicates high uncertainty on the onset and offset of reflectors (both shallow and deeper sections), while showing low uncertainty in the areas of the image with no reflectors.

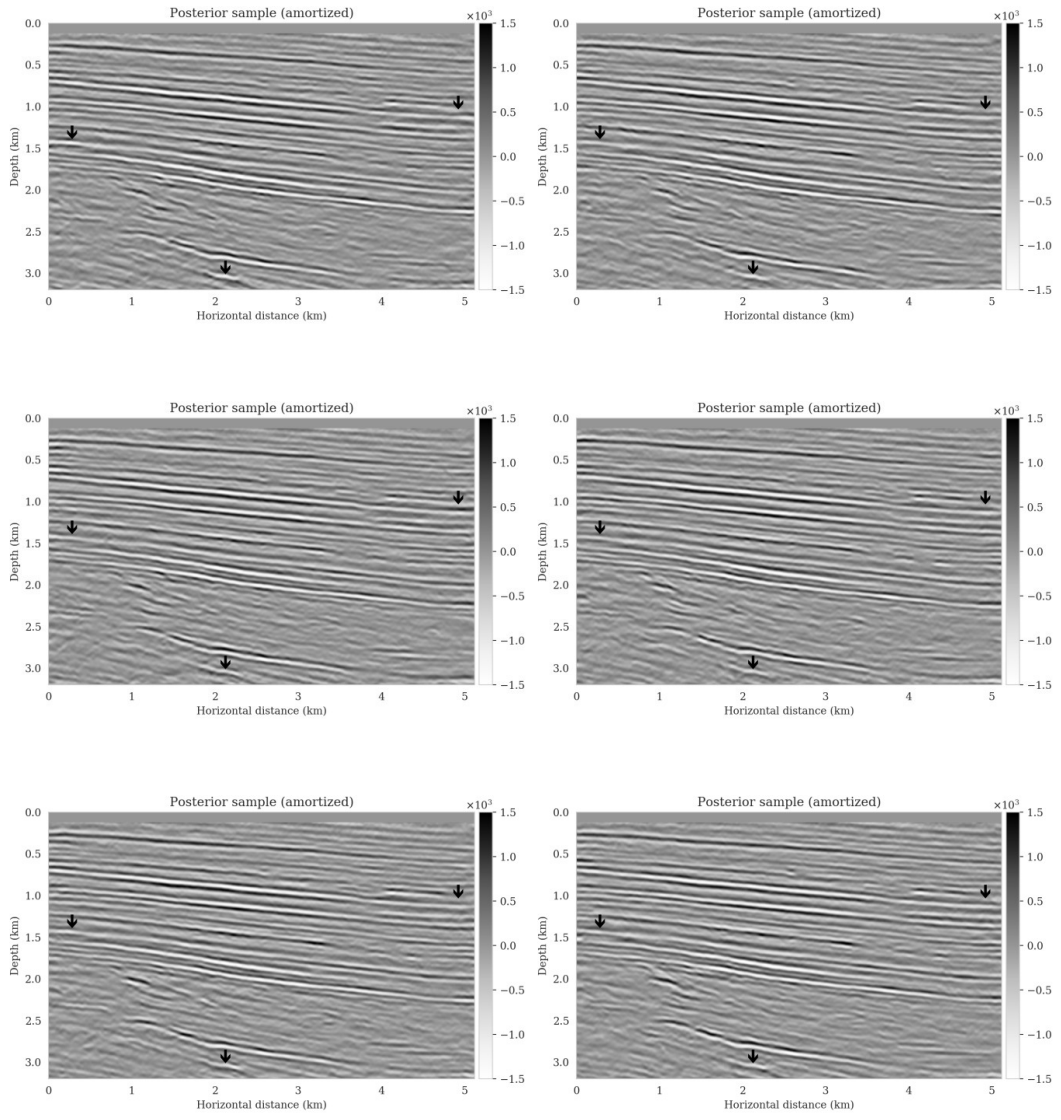
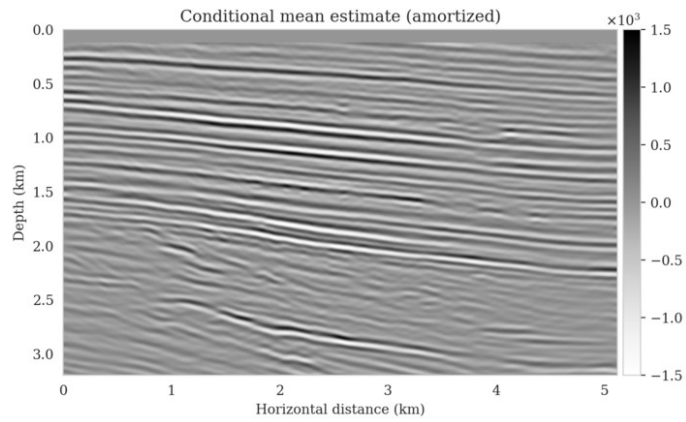
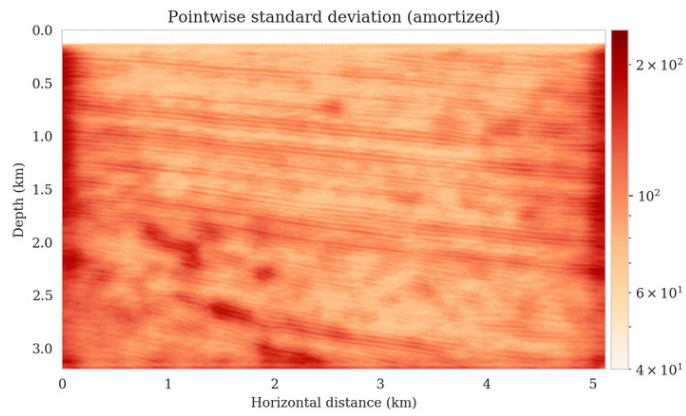


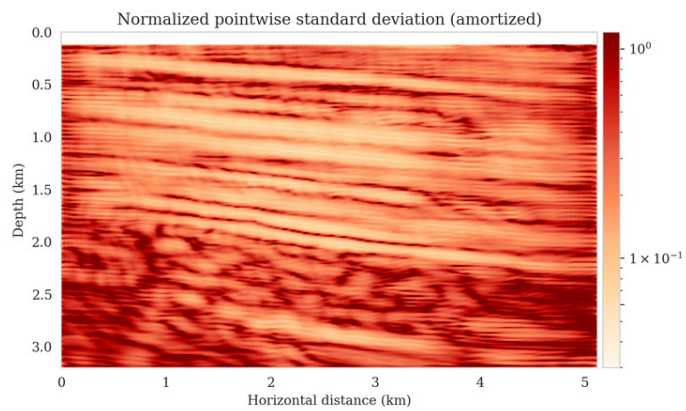
Figure 7: Samples drawn from posterior distribution using the pretrained conditional normalizing flow via equation 9 with SNRs ranging from 8.08 dB to 8.92 dB.



(a)



(b)



(c)

Figure 8: Amortized variational inference results. (a) The conditional (posterior) mean estimate with SNR 9.44 dB. (b) The pointwise standard deviation among samples drawn from the posterior. (c) Normalized pointwise standard deviation by the conditional mean estimate (Figure 8a).

After incurring an upfront cost of training the conditional normalizing flow, the computational cost of sampling the posterior distribution is low as it does not involve any forward operator evaluations. However, the accuracy of the presented results is directly linked to the availability of high-quality training data that fully represent the joint distribution for model and data. Due to our lack of access to the subsurface of the Earth, obtaining high-quality training data is challenging when dealing with geophysical inverse problems. To address this issue, we propose to supplement amortized variational inference with a physics-based latent distribution correction technique that increases the reliability of this approach when dealing with moderate shifts in the data distribution during inference.

## 4 A physics-based treatment to data distribution shifts during inference

For accurate Bayesian inference in the context of amortized variational inference, the surrogate conditional distribution  $p_\phi(\mathbf{x} | \mathbf{y})$  must yield a zero amortized variational inference objective value (equation 8). Achieving this objective is challenging due to lack of access to model and data pairs that sufficiently captures the underlying joint distribution in equation 8. Additionally, due to potential shifts to the joint distribution during inference, i.e., shifts in the prior distribution or the forward (likelihood) model (equation 1), the conditional normalizing flow can no longer reliably provide samples from the posterior distribution due to lack of generalization. Under such conditions feeding latent samples drawn from a standard Gaussian distribution to the conditional normalizing flow may lead to posterior sampling errors. To quantify the posterior distribution approximation error and to propose our correction method, we will use the invariance of the KL divergence to differentiable and invertible mappings [76]. This property relates conditional normalizing flow’s error in posterior distribution approximation to its error in Gaussianizing the input model and data pairs.

### 4.1 KL divergence invariance relation

The errors that the pretrained conditional normalizing flow makes in approximating the posterior distribution can be formally quantified using the invariance of the KL divergence under diffeomorphism mappings [76]. Using this relation, we relate the posterior distribution approximation errors (KL divergence between true and predicted posterior) to the errors that the conditional normalizing flow makes in gaussianizing its inputs (KL divergence between the distribution of “gaussianized” inputs and standard Gaussian distribution). Specifically, for observed data  $\mathbf{y}_{\text{obs}}$  drawn from a shifted data distribution  $\hat{p}_{\text{data}}(\mathbf{y}) \neq p_{\text{data}}(\mathbf{y})$ , the invariance relation states

$$\mathbb{KL}(p_\phi(\mathbf{z} | \mathbf{y}_{\text{obs}}) || \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})) = \mathbb{KL}(p_{\text{post}}(\mathbf{x} | \mathbf{y}_{\text{obs}}) || p_\phi(\mathbf{x} | \mathbf{y}_{\text{obs}})) > 0. \quad (15)$$

In this expression,  $p_\phi(\mathbf{z} | \mathbf{y}_{\text{obs}})$  represents the distribution of conditional normalizing flow output  $\mathbf{z} = f_\phi(\mathbf{x}; \mathbf{y}_{\text{obs}})$ . That is, passing inputs  $\mathbf{x} \sim p(\mathbf{x} | \mathbf{y}_{\text{obs}})$  for one instance of observed data  $\mathbf{y}_{\text{obs}} \sim \hat{p}_{\text{data}}(\mathbf{y})$  to the conditional normalizing flow implicitly defines a (conditional) distribution  $p_\phi(\mathbf{z} | \mathbf{y}_{\text{obs}})$  in the conditional normalizing flow output space. We refer to this distribution as the shifted latent distribution as it is the result of a data distribution shift translated through the conditional normalizing flow to the latent space. The data distribution shifts can be caused by changes in number of sources, noise distribution, wavelet source frequency, and geological features to be imaged. Equation 15 states that the conditional normalizing flow fails to accurately Gaussianize the input models  $\mathbf{x} \sim p(\mathbf{x} | \mathbf{y}_{\text{obs}})$  for the given data  $\mathbf{y}_{\text{obs}}$ . Failure to take into account the mismatch between the shifted latent distribution  $p_\phi(\mathbf{z} | \mathbf{y}_{\text{obs}})$  and  $\mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})$  leads to posterior sampling errors as the KL divergence between the predicted and true posterior distributions is nonzero (equation 15). In other words, feeding latent samples drawn from a standard Gaussian distribution to  $f_\phi^{-1}(\mathbf{z}; \mathbf{y}_{\text{obs}})$  produces samples from  $p_\phi(\mathbf{x} | \mathbf{y}_{\text{obs}})$ , which does not accurately approximate the true posterior distribution under the assumption of data distribution shift. On the other hand, with the same reasoning via the KL divergence invariance relation, feeding samples from the shifted latent distribution  $p_\phi(\mathbf{z} | \mathbf{y}_{\text{obs}})$  to the conditional normalizing flow yields accurate posterior samples. However, obtaining samples from  $p_\phi(\mathbf{z} | \mathbf{y}_{\text{obs}})$  is not trivial as we do not have a closed-form expression for its density. In the next section, we introduce a physics-based approximation to the shifted-latent distribution.

### 4.2 Physics-based latent distribution correction

Ideally, performing accurate posterior sampling via the pretrained conditional normalizing flow—in the presence of data distribution shifts—requires passing samples from the shifted latent distribution



$p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}})$  to  $f_\phi^{-1}(\mathbf{z}; \mathbf{y}_{\text{obs}})$ . Unfortunately, accurately sampling  $p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}})$  requires access to the true posterior distribution, which we are ultimately after and do not have access to. Alternatively, we propose to quantify  $p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}})$  using Bayes’ rule,

$$p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}}) = \frac{p_{\text{like}}(\mathbf{y}_{\text{obs}} \mid \mathbf{z}) p_{\text{prior}}(\mathbf{z})}{\widehat{p}_{\text{data}}(\mathbf{y}_{\text{obs}})}, \quad (16)$$

where the physics-informed likelihood function  $p_{\text{like}}(\mathbf{y}_{\text{obs}} \mid \mathbf{z})$  and the prior distribution  $p_{\text{prior}}(\mathbf{z})$  over the latent variable are defined as

$$\begin{aligned} -\log p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}}) &= -\sum_{i=1}^N \log p_{\text{like}}(\mathbf{y}_{\text{obs},i} \mid \mathbf{z}) - \log p_{\text{prior}}(\mathbf{z}) + \log \widehat{p}_{\text{data}}(\mathbf{y}_{\text{obs}}) \\ &:= \frac{1}{2\sigma^2} \sum_{i=1}^N \|\mathbf{y}_{\text{obs},i} - \mathcal{F}_i \circ f_\phi^{-1}(\mathbf{z}; \mathbf{y}_{\text{obs}})\|_2^2 + \frac{1}{2} \|\mathbf{z}\|_2^2 + \text{const.} \end{aligned} \quad (17)$$

In the above expression, the physics-informed likelihood function  $p_{\text{like}}(\mathbf{y}_{\text{obs}} \mid \mathbf{z})$  follows from the forward model in equation 1 with a Gaussian assumption on the noise with mean zero and covariance matrix  $\sigma^2 \mathbf{I}$ , and the prior distribution  $p_{\text{prior}}(\mathbf{z})$  is chosen as a standard Gaussian distribution with mean zero and covariance matrix  $\mathbf{I}$ . The choice of the likelihood function ensures physics and data fidelity by giving more importance to latent variables that once passed through the pretrained conditional normalizing flow and the forward operator provide smaller data misfits while the prior distribution  $p_{\text{prior}}(\mathbf{z})$  injects our prior beliefs about the latent variable, which is by design chosen to be distributed according to a standard Gaussian distribution.

Due to our choice of the likelihood function and prior distribution above, the effective prior distribution over the unknown  $\mathbf{x}$  is in fact a conditional prior characterized by the pretrained conditional normalizing flow [42]. As observed by Yang and Soatto [77] and Orozco et al. [35], using a conditional prior may be more informative than its unconditional counterpart because it is conditioned by the observed data  $\mathbf{y}_{\text{obs}}$ . Our approach can be also viewed as an instance of online variational Bayes [78] where data arrives sequentially and previous posterior approximates are used as priors for subsequent approximations.

In the next section, we improve the available amortized approximation to the posterior distribution by relaxing the standard Gaussian distribution assumption of the conditional normalizing flow latent distribution.

#### 4.2.1 Gaussian relaxation of the latent distribution

By definition, feeding samples from  $p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}})$  to the pretrained amortized conditional normalizing flows provides samples from the posterior distribution (see discussion beneath equation 15). To maintain the low computational cost of sampling with amortized variational inference, it is imperative that  $p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}})$  is sampled as cheaply as possible. To this end, we exploit the fact that conditional normalizing flows in the context of amortized variational inference are trained to Gaussianize the input model random variable (equation 11). This suggests that the shifted latent distribution  $p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}})$  will be close to a standard Gaussian distribution for a certain class of data distribution shifts. We exploit this property and approximate the shifted latent distribution  $p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}})$  via a Gaussian distribution with an unknown mean and diagonal covariance matrix,

$$p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}}) \approx \mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}, \text{diag}(\mathbf{s})^2), \quad \mathbf{z} \in \mathcal{Z}. \quad (18)$$

In the above expression, the vector  $\boldsymbol{\mu}$  corresponds to the mean and the vector  $\text{diag}(\mathbf{s})^2$  represents a diagonal covariance matrix with diagonal entries  $\mathbf{s} \odot \mathbf{s}$  (with the symbol  $\odot$  denoting elementwise multiplication) that need to be determined. We estimate these quantities by minimizing the reverse KL divergence between the relaxed Gaussian latent distribution  $\mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}, \text{diag}(\mathbf{s})^2)$  and the shifted latent distribution  $p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}})$ . According to the variational inference objective function associated with the reverse KL divergence in equation 5, this correction can be achieved by solving the following

optimization problem (see derivation in Appendix A),

$$\begin{aligned} \boldsymbol{\mu}^*, \mathbf{s}^* &= \arg \min_{\boldsymbol{\mu}, \mathbf{s}} \mathbb{KL} \left( \mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}, \text{diag}(\mathbf{s})^2) \parallel p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}}) \right) \\ &= \arg \min_{\boldsymbol{\mu}, \mathbf{s}} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{z} \mid \mathbf{0}, \mathbf{I})} \left[ \frac{1}{2\sigma^2} \sum_{i=1}^N \left\| \mathbf{y}_{\text{obs}, i} - \mathcal{F}_i \circ f_\phi(\mathbf{s} \odot \mathbf{z} + \boldsymbol{\mu}; \mathbf{y}_{\text{obs}}) \right\|_2^2 \right. \\ &\quad \left. + \frac{1}{2} \left\| \mathbf{s} \odot \mathbf{z} + \boldsymbol{\mu} \right\|_2^2 - \log \left| \det \text{diag}(\mathbf{s}) \right| \right]. \end{aligned} \quad (19)$$

We solve optimization problem 19 with the Adam optimizer where we select random batches of latent variable variables  $\mathbf{z} \sim \mathcal{N}(\mathbf{z} \mid \mathbf{0}, \mathbf{I})$  and data indices. We initialize the optimization problem 19 by  $\boldsymbol{\mu} = \mathbf{0}$  and  $\text{diag}(\mathbf{s})^2 = \mathbf{I}$ . This initialization acts as a warm-start and an implicit regularization [42] since  $f_\phi^{-1}(\mathbf{z}; \mathbf{y}_{\text{obs}})$  for standard Gaussian distributed latent samples  $\mathbf{z}$  provides approximate samples from the posterior distribution—thanks to amortization over different observed data  $\mathbf{y}$ . As a result, we expect the optimization problem 19 to be solved relatively cheaply. Additionally, the imposed standard Gaussian distribution prior on  $\mathbf{s} \odot \mathbf{z} + \boldsymbol{\mu}$  regularizes inversion for the corrections since  $\mathbb{KL}(p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}}) \parallel \mathcal{N}(\mathbf{z} \mid \mathbf{0}, \mathbf{I}))$  is minimized during amortized variational inference (equation 15). To relax the (conditional) prior imposed by the pretrained conditional normalizing flow, instead of a standard Gaussian prior, a Gaussian prior with a larger variance can be imposed on the corrected latent variable. Conditional normalizing flows’ inherent invertibility allows the normalizing flow to represent any solution  $\mathbf{x} \in \mathcal{X}$  in the solution space. This has the additional benefit of limiting the adverse affects of imperfect pretraining of  $f_\phi$  in domains where access to high-fidelity training data is limited, which is often the case in practice. The output of the conditional normalizing flow can be further regularized by including additional regularization terms in equation 19 to prevent it from producing out-of-range, non-physical results. Figure 9 summarizes our proposed method latent distribution correction method.

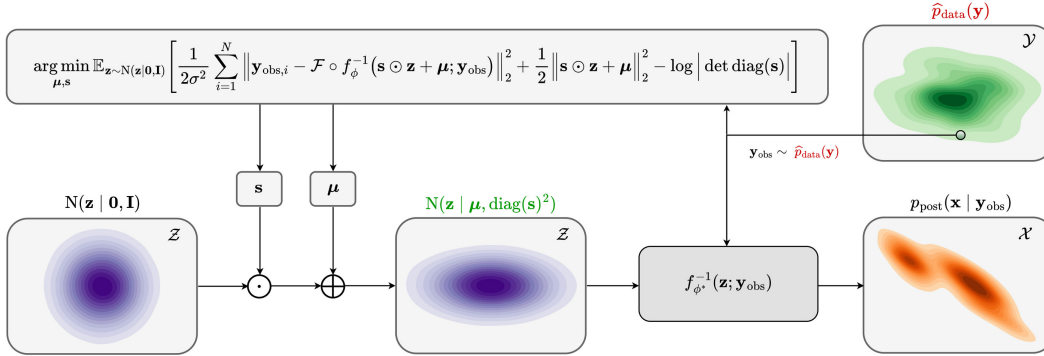


Figure 9: A schematic representation of our proposed method. When dealing with nonzero amortized variational inference objective value (equation 7) or in presence of data distribution shifts during inference, we correct the latent distribution of the pretrained conditional normalizing flow via a diagonal physics-based correction. After the correction, the new latent samples result in corrected posterior samples when fed to the pretrained normalizing flow.

#### 4.2.2 Inference with corrected latent distribution

Once the optimization problem 19 is solved with respect to  $\boldsymbol{\mu}$  and  $\mathbf{s}$ , we obtain corrected posterior samples by passing samples from the corrected latent distribution  $\mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}^*, \text{diag}(\mathbf{s}^*)^2) \approx p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}})$  to the conditional normalizing flow,

$$\mathbf{x} = f_\phi^{-1}(\mathbf{s}^* \odot \mathbf{z} + \boldsymbol{\mu}^*; \mathbf{y}_{\text{obs}}), \quad \mathbf{z} \sim \mathcal{N}(\mathbf{z} \mid \mathbf{0}, \mathbf{I}). \quad (20)$$

These corrected posterior samples are implicitly regularized by the reparameterization with the pretrained conditional normalizing flow and the standard Gaussian distribution prior on  $\mathbf{z}$  [35, 36, 42]. Next section applies this physics-based correction to a seismic imaging example, in which we use the pretrained conditional normalizing flow from the earlier example.

## 5 Latent distribution correction applied to seismic imaging

The purpose of our proposed latent distribution correction approach is to accelerate Bayesian inference while maintaining fidelity to a specific observed dataset and physics. While this method is generic and can be applied to a variety of inverse problems, it is particularly relevant when solving geophysical inverse problems, where the unknown quantity is high dimensional, the forward operator is computationally costly to evaluate, and there is a lack of access to high-quality training data that represents the true heterogeneity of the Earth’s subsurface. Therefore, we apply this approach to seismic imaging to utilize the advantages of generative models for solving inverse problems, including fast conditional sampling and learned prior distributions, while limiting the negative bias induced by shifts in data distributions.

The results are presented for two cases. The first case involves introducing a series of changes in the distribution of observed data, for example changing the number of sources and the noise levels. This is followed by correcting for the error in predictions made by the pretrained conditional normalizing flow using our proposed method. In the second case, in addition to the shifts in the distribution of the observed data (forward model), we also introduce a shift in the prior distribution. We accomplish this by selecting a ground-truth image from a deeper section of the Parihaka dataset that has different image characteristics than the training images, such as tortuous reflectors and more complex geological features. In both cases, we expect the outcome of the Bayesian algorithm to improve following the correction of latent distributions described above. We provide qualitative and quantitative evaluations of the Bayesian inference results.

### 5.1 Shift in the forward model

In the following example, we introduce shifts in the distribution of observed data—compared to the pretraining phase—by changing the forward model. The shift involves reducing the number of sources ( $N$  in equation 1) by a factor of two to four, while adding band-limited noise with 1.5 to three times larger standard deviation ( $\sigma$  in equation 3). We will demonstrate the potential pitfalls of relying solely on the pretrained conditional normalizing flows in circumstances where the distribution of observed data has shifted. With the use of our latent distribution correction, we will demonstrate that we are able to correct for errors that are made by the pretrained conditional normalizing flow as a result of changes to data distribution.

Following the description of the problem setup, we will also provide comparisons between the conditional mean estimation quality before and after the latent distribution correction step. Before moving on to our results relating to uncertainty quantification on the image, we demonstrate the importance of the correction step by visualizing the improvements in fitting the observed data. Lastly, we perform a series of experiments to verify our Bayesian inference results.

#### 5.1.1 Problem setup

To induce shifts in the data distribution, we reduce the number of sources and increase the standard deviation of the added band-limited noise. Consequently, we have reduced the amount of data (due to having fewer source experiments) and decreased the SNR of each shot record (due to being contaminated with stronger noise). As a consequence, seismic imaging becomes more challenging, i.e., more difficult to estimate the ground truth image, and it is expected that the uncertainty associated with the problem will also increase.

We use the same ground-truth image as in the previous example (Figure 6a), while experimenting with 25, 51, 102 sources and adding band-limited noise that has 1.5, 2.0, 2.5, and 3.0 times larger standard deviation than the pretraining setup. For each combination of source number and noise level (12 combinations in total), we compute the reverse-time migrated image corresponding to that combination. Next, we perform latent distribution corrections for each of the 12 seismic imaging instances. All latent distribution correction optimization problems (equation 19) are solved using the Adam optimization algorithm [62] for five passes over the shot records (epochs). We did not observe a significant decrease in the objective function after five epochs. The objective function is evaluated each iteration by drawing a single latent sample from the standard Gaussian distribution and randomly selecting (without replacement) a data index  $i \in \{1, \dots, 25\}$ . We use a stepsize of  $10^{-1}$  and decrease it by a factor of 0.9 at the end of every two epochs.

After solving the optimization problem 19 for the different seismic imaging instances, we obtain corrected posterior samples for each instance. The next section provides a detailed discussion of the latent distribution correction that was applied to one such instance that had a significant shift in data distribution.

### 5.1.2 Improved Bayesian inference via latent distribution correction

The aim of this section is to demonstrate how latent distribution correction can be used to mitigate errors induced by data distribution shifts. Specifically, we present the results for the case where the number of sources is reduced by a factor of four ( $N = 25$ ) as compared to the pretrained data generation setup. The 25 sources are spread periodically over the survey area with a source sampling of approximately 200 meters. Moreover, we contaminate the resulting shot records with band-limited noise with an increased standard deviation of 2.5 times when compared to the pretraining phase. The overall SNR for the data thus becomes  $-2.78$  dB, which is 7.95 dB lower than the SNR of the observed data during pretraining (Figure 5b). Figure 10 shows one of the ‘25 shot records.

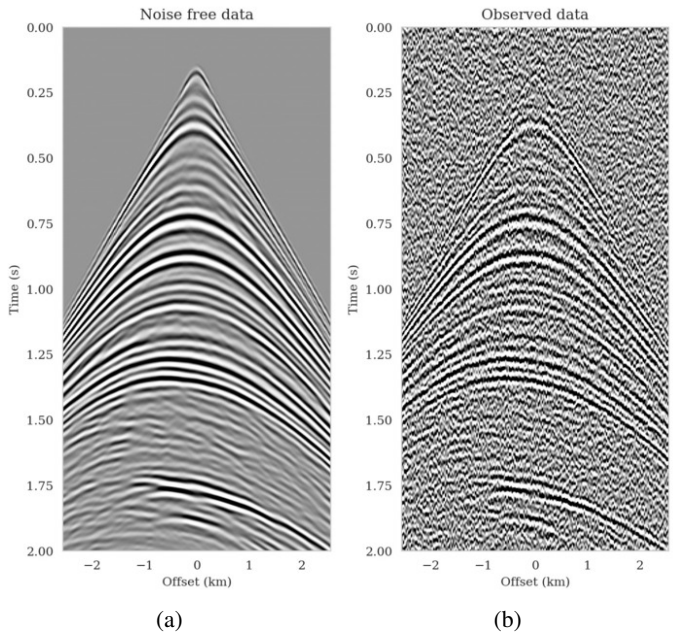
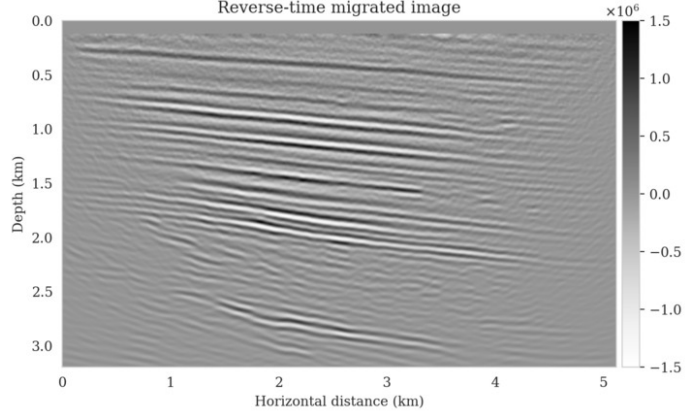


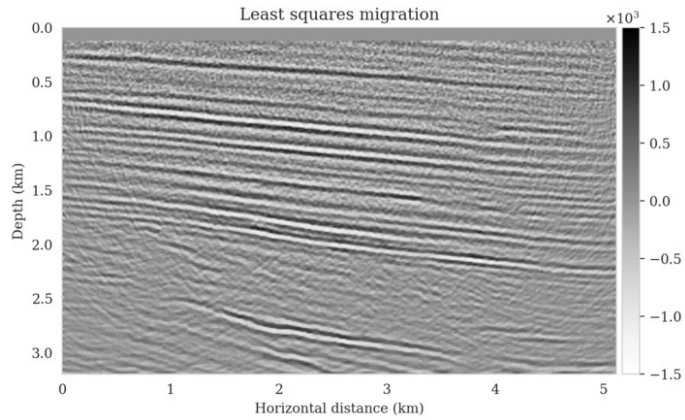
Figure 10: A shot record from the shifted data distribution. (a) Noise-free linearized data (same as Figure 5a). (b) Noisy linearized data with 2.5 larger band-limited noise standard deviation (SNR  $-2.78$  dB).

Utilizing the above observed dataset, we compute the reverse-time migrated image as an input to our pretrained conditional normalizing flow (Figure 11a). In contrast to the reverse-time migrated image shown in Figure 5b, this migrated image is, as expected, noisier, and it displays visible near-source imaging artifacts as a result of coarse source sampling. Additionally, we compute the least-squares migrated seismic image that is obtained by minimizing the negative-log likelihood (see the likelihood term in equation 3). This image, shown in Figure 11b, was constructed by fitting the data without incorporating any prior information. It is evident from this image that there are strong artifacts caused by noise in the data, underscoring the importance of incorporating prior knowledge into solving seismic imaging.

**5.1.2.1 Improvements in posterior samples and conditional mean estimate** To obtain amortized (uncorrected) posterior samples, we feed the reverse-time migrated image (Figure 11a) and latent samples drawn from the standard Gaussian distribution to the pretrained normalizing flow. These samples, which are shown in the left column of Figure 12, contain artifacts near the top of the image. These artifacts are related to the near-source reverse-time migrated image artifacts (Figure 11a). Since the reverse-time migrated images used during pretraining do not contain near-source



(a)



(b)

Figure 11: Latent distribution correction experiment setup. (a) Reverse-time migrated image corresponding to the shifted forward model with SNR  $-8.22$  dB. (b) Least squares imaging, which is equivalent to the minimizing  $\sum_{i=1}^N \|\mathbf{d}_i - \mathbf{J}(\mathbf{m}_0, \mathbf{q}_i)\delta\mathbf{m}\|_2^2$  with respect to  $\delta\mathbf{m}$  with no regularization. The SNR for this estimate is  $6.90$  dB.

imaging artifacts—due to fine source sampling—the pretrained normalizing flow fails to eliminate them. Further, the uncorrected posterior samples do not accurately predict reflectors as they approach the boundaries and deeper sections of the image.

To illustrate the improved posterior sample quality following latent distribution correction, we feed latent samples drawn from the corrected latent distribution to the pretrained normalizing flow (right column of Figure 12). Comparing the left and right columns in Figure 12 indicates an improvement in the quality of samples from the posterior distribution, which can be attributed to the attenuation of near-top artifacts and an improvement in the image quality close to the boundary and deeper reflectors in the image. Moreover, the SNR values of the posterior samples after correction are approximately 3 dB higher, which represents a significant improvement.

To compute the conditional mean estimate, we simulate one thousand posterior samples before and after latent distribution correction. As with the posterior samples before correction, drawing samples after correction is very cheap once the correction is done as it only requires evaluating the conditional normalizing flow over the corrected latent samples. Figures 13a and 13b show conditional mean estimates before and after latent distribution correction, respectively. The conditional mean estimate before correction reveals similar artifacts as the posterior samples before correction, in particular, near-top imaging artifacts due to coarse sources sampling and less illumination of reflectors located closer to the boundary and deeper portions of the image. The importance of our proposed

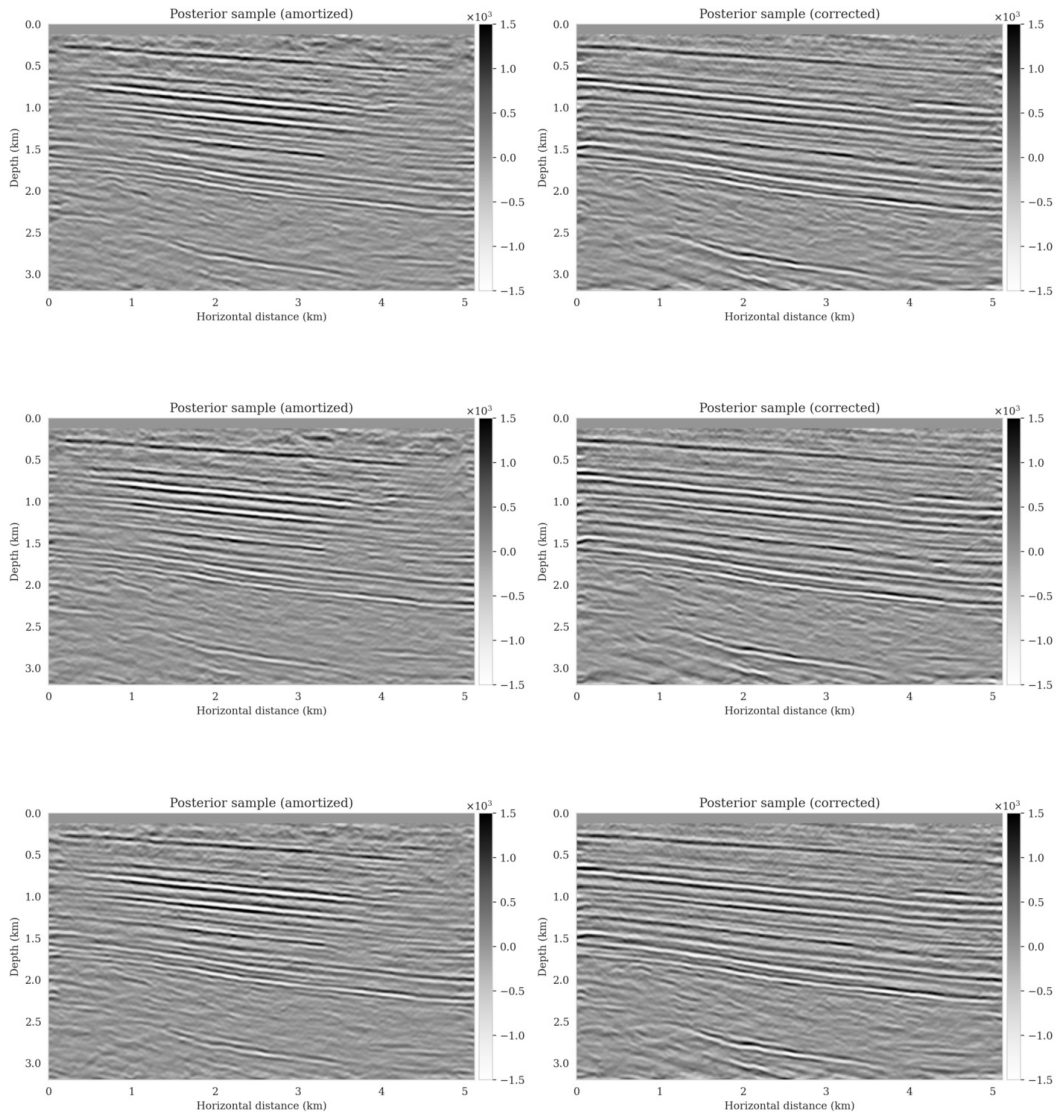


Figure 12: Samples from the posterior distribution (left) without latent distribution correction with SNRs ranging from 4.57 dB to 5.21 dB; and (right) after latent distribution correction with SNRs ranging from 7.80 dB to 8.53 dB.

latent distribution correction can be observed by juxtaposing the conditional mean estimate before (Figure 13a) and after correction (Figure 13b). The conditional mean estimate obtained after latent distribution correction eliminates the aforementioned inaccuracies and enhances the quality of the image by approximately 4 dB. We gain similar improvements in SNR compared to the least-squares migrated image (Figure 11b) with virtually the same cost, i.e., five passes over the shot records. This significant improvement in SNR also is complimented by access to information regarding the uncertainty of the image.

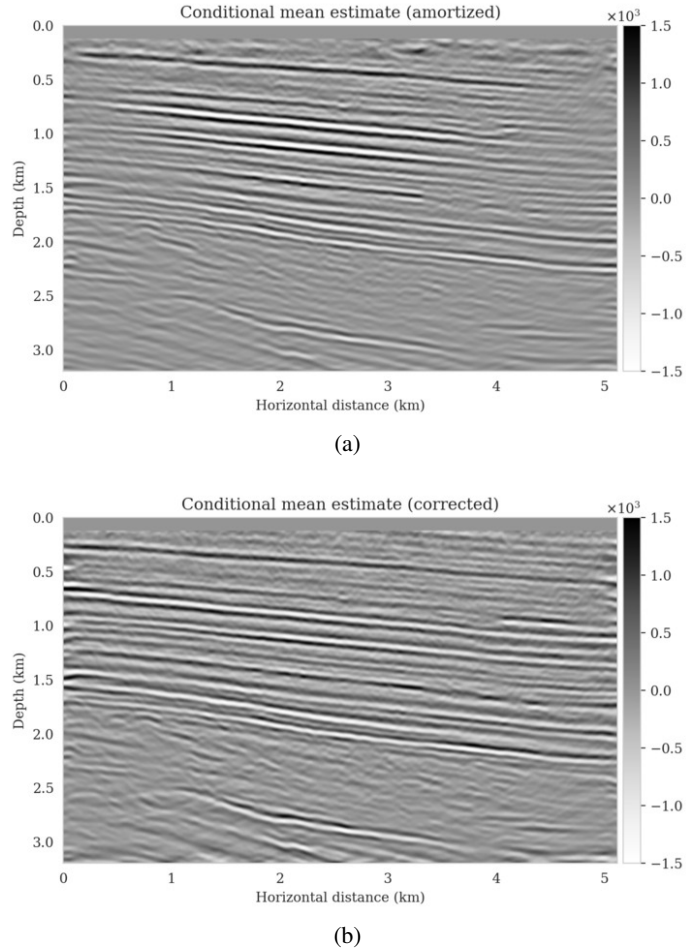


Figure 13: Improvements in conditional mean estimate due to latent distribution correction. (a) The conditional (posterior) mean estimate using the pretrained conditional normalizing flow without correction (SNR 6.29 dB). (b) The conditional mean estimate after latent distribution correction (SNR 10.36 dB).

**5.1.2.2 Data-space quality control** As the latent distribution correction step involves finding latent samples that are better suited to fit the data (equation 19), we can expect an improvement in fitting the observed data after correction. Predicted data is obtained by applying the forward operator to the conditional mean estimates, before and after latent distribution correction. Figures 14a and 14b show the predicted shot records before and after correction, respectively. In spite of the fact that both predicted data appear to be similar to ideal noise-free data (Figure 10a), the data residual associated with the conditional mean without correction reveals several coherent events that contain valuable information about the unknown seismic image. The latent distribution correction allows us to fit these coherent events as indicated by the data residual associated with the corrected conditional mean estimate.

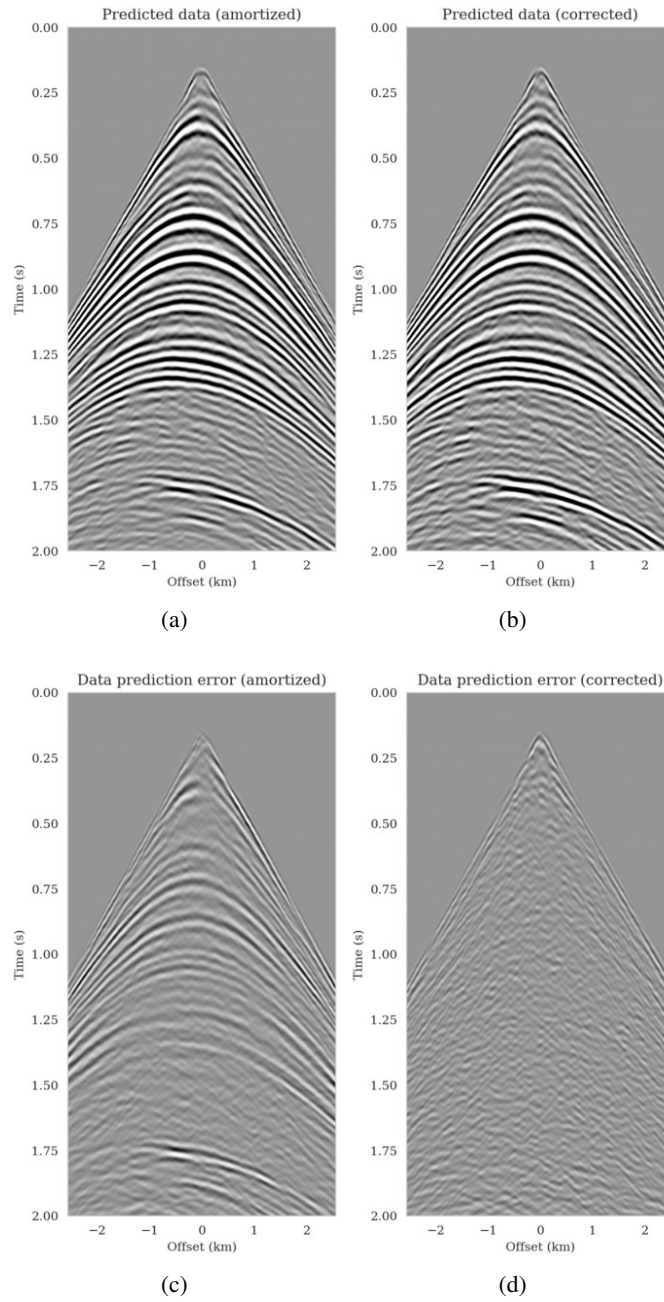


Figure 14: Quality control in data space. Data is simulated by applying the forward operator to the conditional mean estimate (a) before (SNR 11.62 dB); and (b) after latent distribution correction (SNR 16.57 dB). (c) Prediction errors associated with Figure 14a. (d) Prediction errors associated with Figure 14b (after latent distribution correction).



**5.1.2.3 Uncertainty quantification—pointwise standard deviation and histograms** We exploit cheap access to corrected samples from the posterior in order to extract information regarding uncertainty in the image estimates. Figure 15a displays the pointwise standard deviation among the one thousand corrected posterior samples. The overprint by the strong reflectors can be reduced by normalizing the standard deviation using a stabilized division by the conditional mean (Figure 15b). The pointwise standard deviation plots indicate high uncertainty in areas near the boundaries of the image and in the deep parts of the image where illumination is relatively poor. This observation is more evident in Figure 16, which displays three vertical profiles as 99% confidence intervals (orange colored shading) illustrating the expected increasing trend of uncertainty with depth. We additionally observe that the ground truth (dashed black) falls within the confidence intervals for most of the areas.

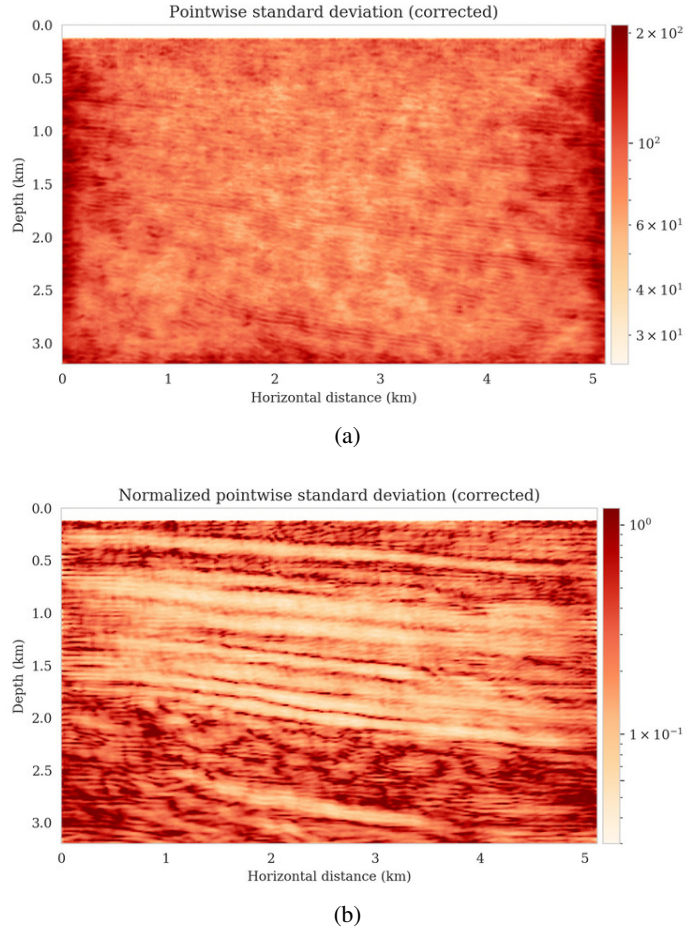


Figure 15: Uncertainty quantification with latent distribution correction. (a) The pointwise standard deviation among samples drawn from the posterior after latent distribution correction. (b) Normalized pointwise standard deviation by the conditional mean estimate (Figure 13b).

To demonstrate how the corrected posterior is informed by the observed data, we calculated histograms at three locations in Figure 15a. Prior histograms are calculated by feeding latent samples drawn from the standard Gaussian distribution to the pretrained conditional normalizing flow without using data conditioning (see Kruse et al. [28] and Siahkoohi and Herrmann [33] for more information). These samples in the image spaces are indicative of samples from the prior distribution implicitly learned by the conditional normalizing flow during pretraining. The resulting prior histograms are shown in Figure 17. Corresponding histograms are also obtained for the uncorrected amortized posterior distribution (equation 12). As mentioned before, the uncorrected posterior distribution serves as an implicit conditional prior for the subsequent step of correction of the latent distribution. The green histograms in Figure 17 represent the uncorrected amortized posterior distribution. A similar procedure is followed to obtain histograms after latent distribution correction (blue histograms). As

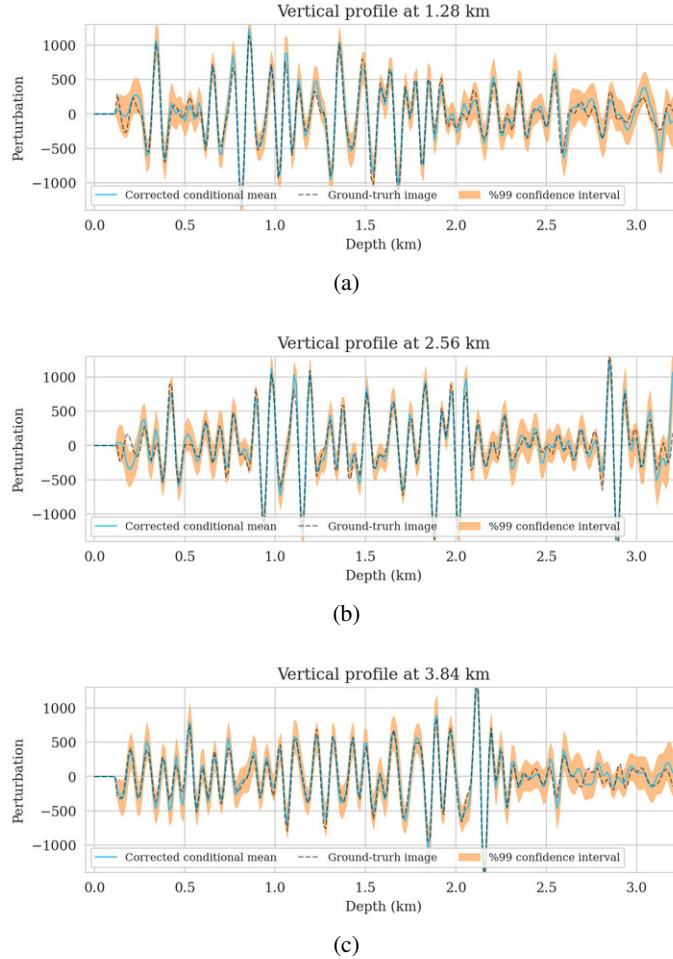
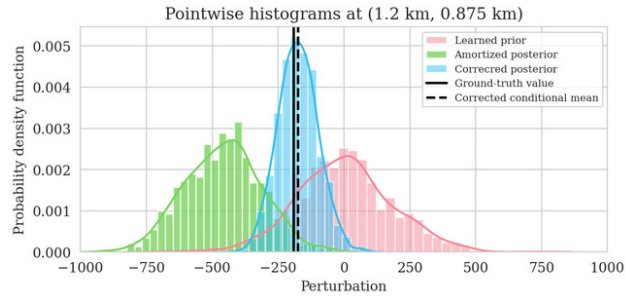


Figure 16: Confidence intervals for three vertical profiles. Traces of 99% confidence interval (shaded orange color), corrected conditional mean (solid blue), and ground truth (dashed black) at (a) 1.28 km, (b) 2.56 km, and (c) 3.84 km horizontal location.

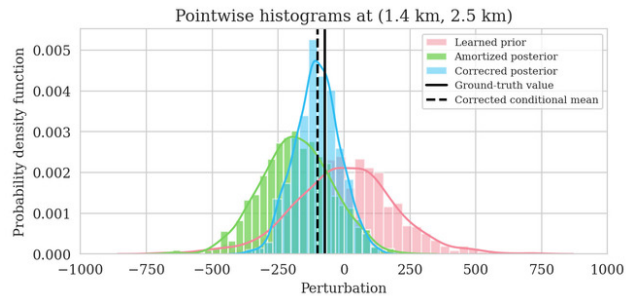
expected, the histograms of the posterior distribution are considerably narrower than those of the learned prior, which indicates that the posterior is further informed by the specific observed dataset and physics. As a means of evaluating the effect of latent distribution correction, we provide a vertical solid line showing the ground truth value's location. All three corrected posterior histograms for each location are shifted towards the ground truth, and their (conditional) mean plotted with the dashed vertical line indicates improved recovery of the ground truth. Compared to the amortized uncorrected histograms, the corrected histograms in Figures 17a – 17b are further contracted, suggesting that the latent distribution correction step has further informed the inference by the data.

### 5.1.3 Bayesian inference verification

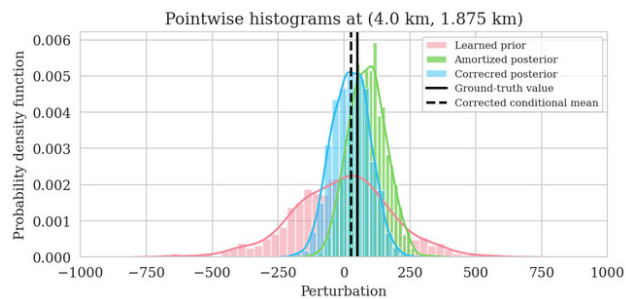
While we investigated the accuracy of the conditional mean estimate after correction, we do not have access to the underlying true posterior distribution to verify our proposed posterior sampling method. This is partly due to our learned prior and the implicit conditional prior used in latent variable correction, which make traditional MCMC-based comparisons challenging. To further validate our Bayesian inference procedure, we conduct a series of experiments in which we investigate the effect of gradual increase in the number of sources ( $N$ ) and reduction of the noise level. As the number of sources increases and the noise level decreases, we expect to see an increase in seismic image quality and a decrease in uncertainty.



(a)



(b)



(c)

Figure 17: Pointwise prior (red), uncorrected amortized posterior (green), and latent distribution corrected posterior (blue) histograms along with the true perturbation values (solid black line) and the corrected conditional mean (dashed black line) for points located at (a) (1.2 km, 0.875 km), (b) (1.4 km, 2.5 km), and (c) (4.0 km, 1.875 km).

**5.1.3.1 Estimation accuracy** The accuracy of the Bayesian parameter estimation method is directly affected by the amount of data that has been collected [2]. That is to say with more observed data (larger  $N$ ), we should be able to obtain a more accurate seismic image estimate. The same principle allows us to stack out noise when increasing the fold in seismic data acquisition. In order to assess whether our Bayesian inference approach has this property, we repeat the latent distribution correction process while varying the number of sources (using  $N = 25, 51, 102$  sources) and the amount of band-limited noise (standard deviations 1.5, 2.0, 2.5, and 3.0 times greater than the noise standard deviation during pretraining). Each of the 12 instances of latent distribution correction problems is treated similarly with respect to the number of passes made over the shot records and other optimization parameters. For each of the 12 combinations of source numbers and noise levels, we calculate the corrected conditional mean estimate and plot the SNRs as a function of the noise’s standard deviation in Figure 18.

For a fixed number of sources (25, 51, and 102 sources shown with red, green, and blue colors respectively), we plot the corrected conditional means SNR as a function of the noise standard deviation. There is a clear increase in SNR trend as we decrease the noise level. In the same way, for each fixed noise level, the SNR increases with the number of sources. This verifies the our Bayesian inference method yields a more accurate estimate of the conditional mean for larger number of sources and smaller noise levels.

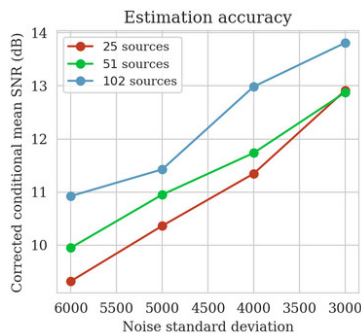


Figure 18: Estimation accuracy as a function of number of sources and noise levels. Colors correspond to different source numbers.

**5.1.3.2 Bayesian posterior contraction** An alternative Bayesian inference verification method involves analyzing the Bayesian posterior contraction, that is, the decrease of uncertainty with more data. To examine whether or not our Bayesian inference method possesses this property, we visually inspect the resulting pointwise standard deviation plots in Figure 19 for the 12 possible combinations of source numbers and noise levels. Each row corresponds to the pointwise standard deviation plot for a fixed noise standard deviation ( $\sigma$ ), where the number of sources ( $N$ ) decreases from left to right. In each column, we maintain the number of sources and we plot the pointwise standard deviation as we increase the noise standard deviation from top to bottom. There is a consistent increase in standard deviation values as we move from the top-left to the bottom-right corner. In other words, the posterior contract (shrinks) when we have more data (more numbers of sources, Figure 19 from right to left) and when we have less noise (Figure 19 from bottom to top), which effectively means more data.

Figure 20 offers an alternate method of visualizing posterior contraction, displaying box plots of the standard deviation values for each of the 12 images in Figure 19. In each of the three box plots, the vertical axis corresponds to the noise standard deviation, and the horizontal axis represents the possible values in the posterior pointwise standard deviation plots. The box indicates the values that are between the first and third quartiles (where half of the possible values fall) and the line in the middle indicates the median value. Figures 20a to 20c show box plots for experiments with 25, 51, and 102 sources, respectively, with each box plot color reflecting a particular noise level. In each of the Figures 20a to 20c, we observe a decrease in the range of posterior standard deviation values, including median and quartiles, as we lower the noise level from left to right. Similarly, for the same noise levels, that is, box plots of the same color, the standard deviations decrease from Figure 20a to

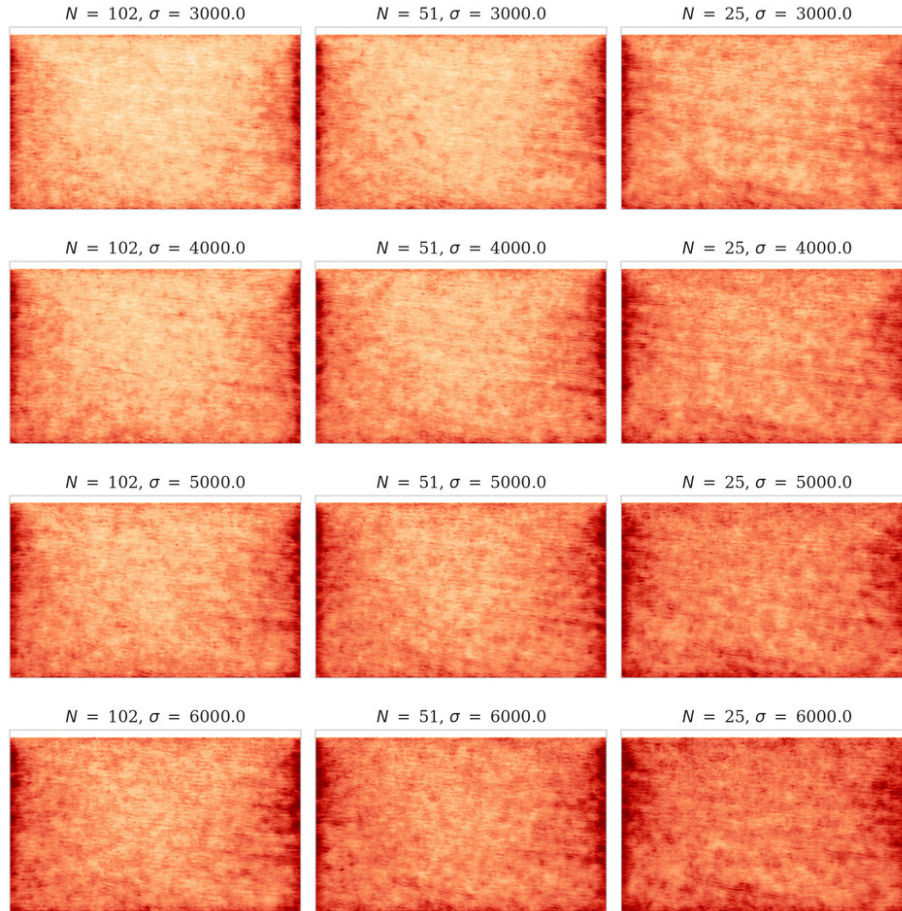


Figure 19: Bayesian posterior contraction: visual inspection. Pointwise standard deviations for varying number of sources (decreasing left to right) and noise variances (increases top to bottom).

Figure 20c (increasing number of sources). The observed trends in Figures 19 and 20 verify that our Bayesian inference method exhibits the Bayesian posterior contraction property.

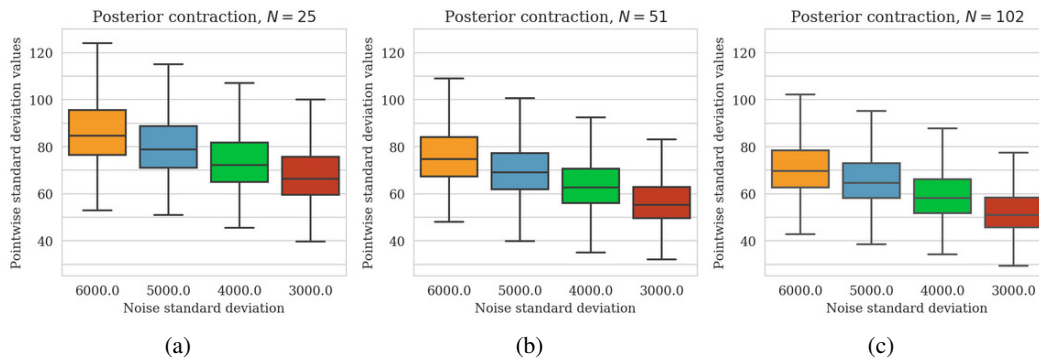


Figure 20: Box plots of pointwise standard deviation values as a function of noise level for number of sources (a)  $N = 25$ , (b)  $N = 51$ , and (c)  $N = 102$ .

## 5.2 Shift in the forward model and prior distribution

This example is intended to demonstrate how our latent distribution correction method can perform Bayesian inference when the unknown ground truth image has properties that differ from the seismic images in the pretraining dataset. This expectation is based on the invertible nature of conditional normalizing flows, which allows them to represent any image in the image space [42].

As a means to mimic this scenario, unlike images used in pretraining, we have extracted two 2D seismic images from deeper sections of the Parihaka dataset. In these sections, there are fewer continuous reflectors and more complex geological features. Figures 21a and 21b show two images with drastically different geological features when compared to Figure 6a. Following the pretraining data acquisition setup, we add a water column on top of these images to reduce the near-source artifacts. Similar to the previous experiment involving forward model shifts, we use only 25 sources with 200 meters sampling distance and add noise with a 2.5 fold larger standard deviation than during the pretraining phase. With this setup, the noisy observed data associated with experiments involving seismic images in Figures 21a and 21b have a SNR of  $-1.56$  dB and  $-2.41$  dB, respectively.

As part of our analysis, we compute the reverse-time (second row in Figure 21) and least-squares (third row in Figure 21) migrated images for these two ground-truth images, where the former images serve as inputs to the pretrained conditional normalizing flows. Similar to the previous example, the reverse-time migrated images contain the near-source artifacts and are contaminated by the input noise. Moreover, the least-squares migrated images highlight the importance of including prior information in this imaging problem, since this image contains strong noise-related artifacts, which might impact downstream tasks, such as horizon tracking.

### 5.2.1 Conditional mean and pointwise standard deviations

To obtain samples from the posterior distribution, we feed the reverse-time migrated images (Figures 21c and 21d) into the pretrained conditional normalizing flow. The posterior is sampled using either standard Gaussian distributions or corrected latent samples. It is apparent, once more, that the uncorrected conditional mean estimates are significantly contaminated by artifacts in the near-source region (Figures 22a and 22b). Another type of noticeable error in these predicted images includes lower amplitudes in deeper and closer to boundary reflectors. A comparison of the conditional means estimates before (Figures 22a and 22b) and after (Figures 22c and 22d) latent distribution correction indicates attenuation of near-source artifacts as well as an improvement in reflector illumination near the boundary and deeper sections where images are more difficult to capture. Following latent distribution correction, the conditional mean estimate SNR is improved by three to four decibels for both images. To provide more quantitative results, we ran the experiments as set up in this section for eight additional seismic images sampled from deeper sections of the Parihaka dataset, sampled from deeper sections of the Parihaka dataset. The SNR of the estimated seismic images before and after latent distribution are  $5.91 \pm 0.49$  dB and  $9.15 \pm 0.73$  dB, respectively.

Careful inspection at the boundaries in the corrected conditional mean estimates reveals some nonrealistic events near the boundaries. The plots of pointwise standard deviations (Figures 23a and 23b) associated with corrected conditional means, however, clearly indicate that there is uncertainty for these events. This illustrates the importance of uncertainty quantification and not relying on a single estimate when addressing ill-posed inverse problems. To diminish the imprint of strong reflectors in the pointwise standard deviations plots, we also display these images when normalized with respect to the envelopes of the conditional mean estimates in Figures 23c and 23d.

### 5.2.2 Data residuals

As before, we confirm that the latent distribution correction improves the fit of the data. Figure 24 shows the predicted data as well as the data residuals for all conditional mean estimates. The predicted data are obtained by applying the forward operator to the conditional mean estimates, both before and after latent distribution correction. The predicted data before and after correction are presented in the first and second columns, respectively. The corresponding data residuals before and after correction, which are computed by subtracting the predicted data from ideal noise-free data, can be seen in the third and last columns of Figure 24. Evidently, the latent distribution correction stage has resulted in a better fit with the observed data, as coherent data events show up in Figures 24c and 24g, but are attenuated in the corrected residual plots (Figures 24d and 24h).

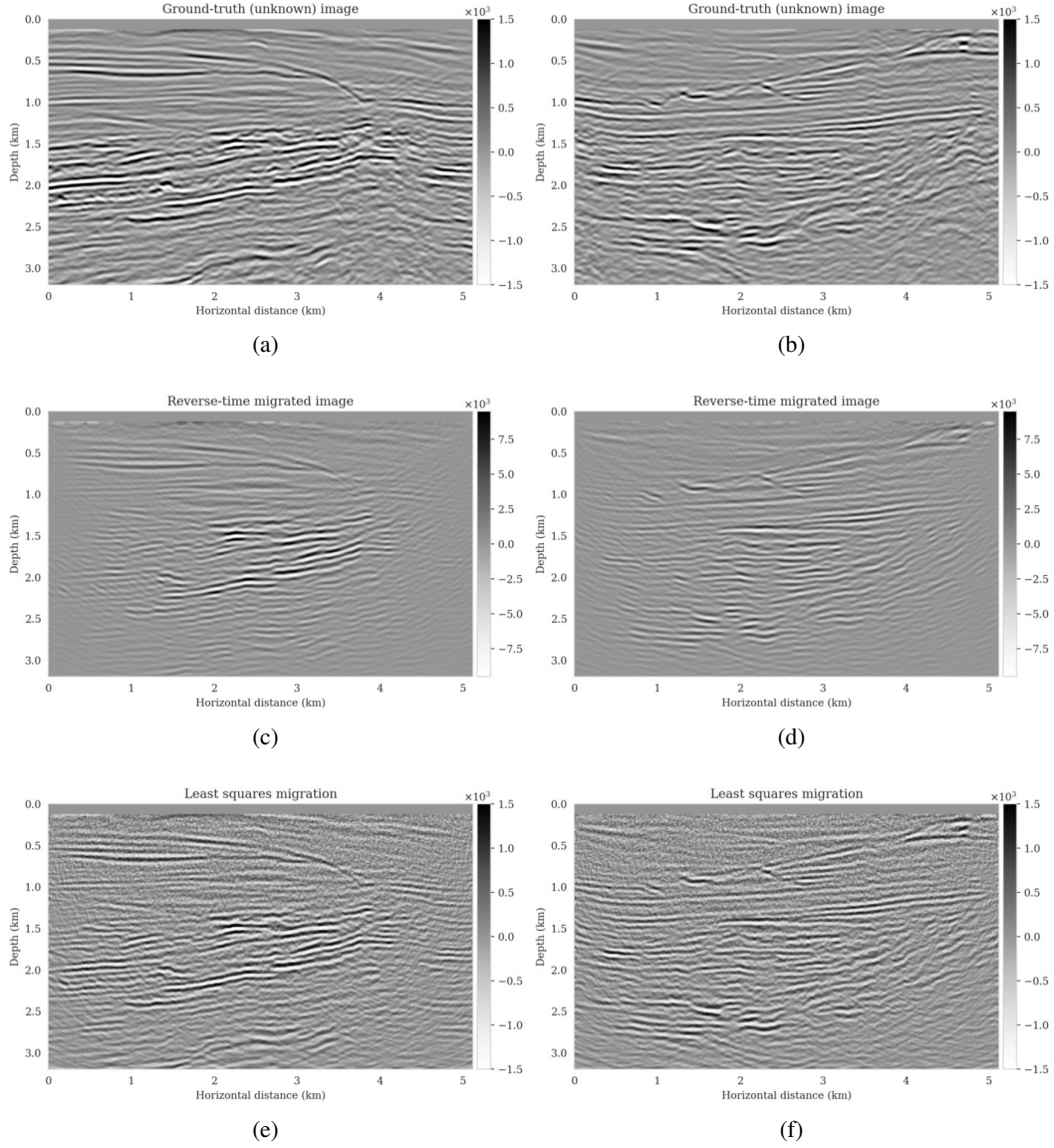


Figure 21: Setup for experiments involving shifts in the prior distribution. (a) and (b) Two high-fidelity ground-truth images from deeper sections of the Parihaka dataset. (c) and (d) Reverse-time migrated image corresponding ground-truth images in Figures 21a (SNR  $-8.02$  dB) and 21b (SNR  $-8.49$  dB), respectively. (e) and (f) Least squares imaging results (no regularization) corresponding ground-truth images in Figures 21a (SNR  $4.94$  dB) and 21b (SNR  $5.59$  dB), respectively.

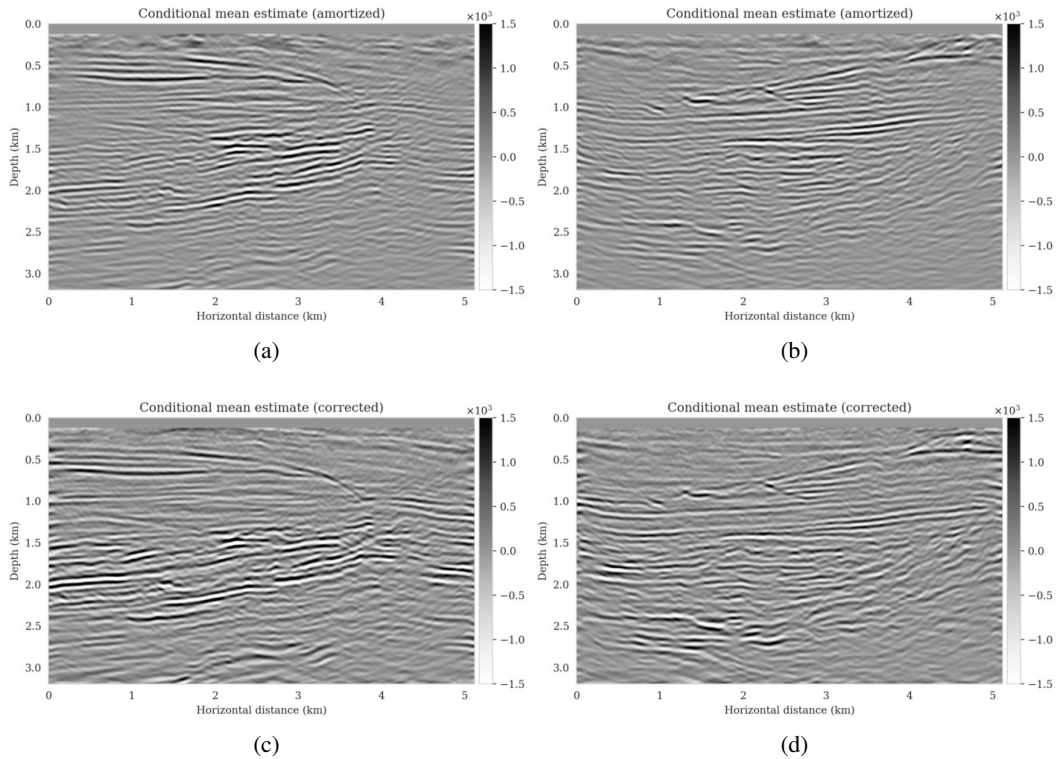


Figure 22: The improvement in conditional mean estimate due to latent distribution correction. (a) and (b) Amortized uncorrected conditional (posterior) mean estimates with SNRs 5.47 dB and 6.17 dB, respectively. (c) and (d) Conditional (posterior) mean estimates after latent distribution correction with SNRs 9.40 dB and 9.11 dB, respectively.

## 6 Discussion

The examples presented demonstrate that deep neural networks trained in the context of amortized variational inference can facilitate solving inverse problems in two ways: (1) incorporating prior knowledge gained through pretraining; and (2) accelerating Bayesian inference and uncertainty quantification. Despite the fact that amortized variational inference is capable of sampling the posterior distribution without requiring forward operator evaluations during inference, the extent to which it is considered reliable is dependent upon the availability of high quality training data. We demonstrate this limitation of amortized variational inference via a seismic imaging example where we alter the number of sources, variance of noise, and the present geological features to be imaged in comparison to the pretraining phase. These variations shift the data distribution, and in certain cases, e.g., imaging slightly different geological features, this can be also thought as evaluating the pretrained conditional normalizing flow over data samples from the low probability regions of the training distribution. Due to shifts in the data distribution, the obtained posterior samples via amortized variational inference in our numerical experiments had unusual near-source artifacts and less illuminated reflectors, which we can be partly attributed to lack of these artifacts in the training reverse-time migrated images.

As part of our efforts to extend the application of the supervised amortized variational inference methods to domains with limited access to high-quality training data, we developed an unsupervised, physics-based variational inference formulation over the latent space of a pretrained conditional normalizing flow that mitigates some of the posterior sampling errors induced by data distribution shifts. In this approach, a diagonal correction to the latent distribution is learned that ensures that the conditional normalizing flow output distribution better matches the desired posterior distribution. Based on observations and other research [35, 36, 42], we found that normalizing flows, because of their invertibility, are capable of partially mitigating errors related to changes in data distributions



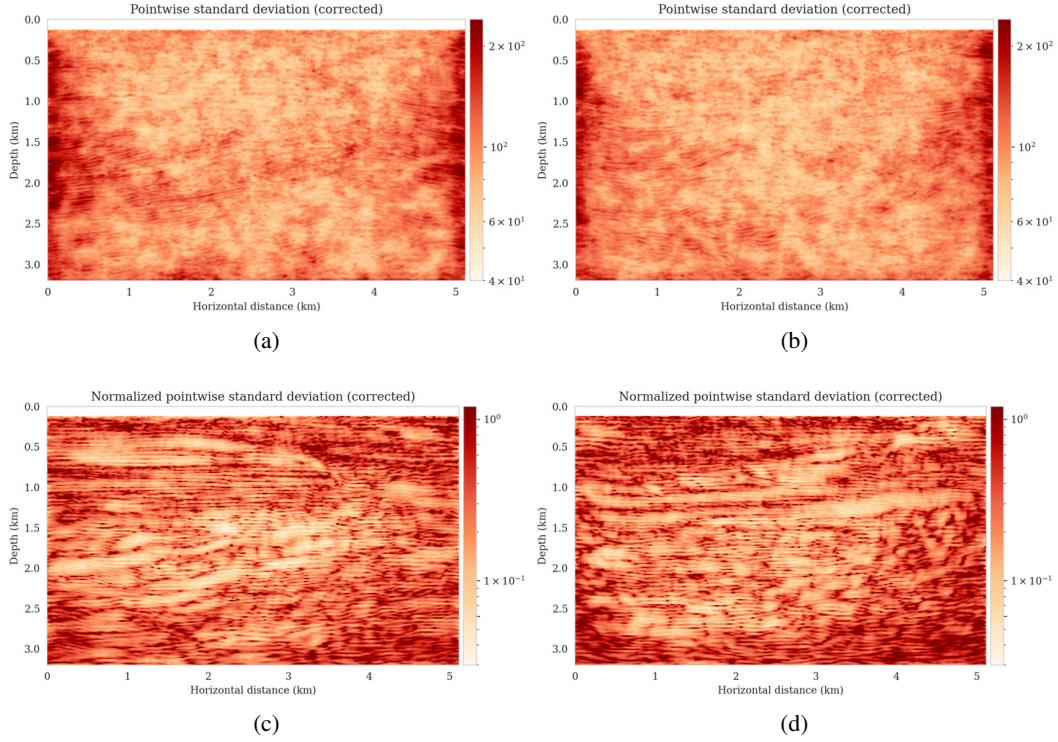


Figure 23: Uncertainty quantification with latent distribution correction. (a) and (b) The pointwise standard deviation estimates among samples drawn from the posterior after latent distribution correction. (b) Pointwise standard deviations normalized by the envelop of conditional mean estimates.

when they are used to solve inverse problems. We leave the delineation of which data distribution shifts can be handled by our diagonal correction approach to future work. Needless to say, by adhering to more complex transformations in the latent space, e.g., via a neural network, a wide range of data distribution shifts can be handled but it would potentially require a more computationally costly correction step.

Latent distribution correction requires several passes over the data (five in our example), which includes evaluation of the forward operator and the adjoint of its Jacobian. Due to the amortized nature of our approach, these costs are significantly lower than those incurred by other non-amortized variational inference methods [19, 46, 48, 49], specifically tens of thousands of forward operator evaluations reported by Zhao et al. [47] in the context of travel-time tomography. By amortizing over the data, the pretrained conditional normalizing flow provides posterior samples for previously unseen data (drawn from the same distribution as training data) without the need for latent distribution correction. In presence of moderate data distribution shifts, the learned posterior is adjusted to new observed data via the latent distribution correction step, which could be considered an instance of transfer learning [79]. In contrast, existing methods that perform variational inference in the latent space [50, 52] require a more substantial modification to the latent distribution, as the pretrained model in these methods originally provides samples from the prior distribution. For large-scale inverse problems, such as seismic imaging, where solving a partial differential equation is required to evaluate the forward operator, reduced computational costs of Bayesian inference are particularly important. Our approach reduces the computational costs associated with Bayesian inference in such problems while also complementing the inversion with a learned prior. In order to quantify the extent to which a diagonal correction to the latent distribution mitigates errors resulting from data distribution shifts, more research is required.

While the examples we presented in this paper were in regards to linear inverse problem, variational inference can be also applied to nonlinear problems [20, 23, 25, 48, 80]. However, in the context of amortized variational inference there are two main considerations that need to be taken into

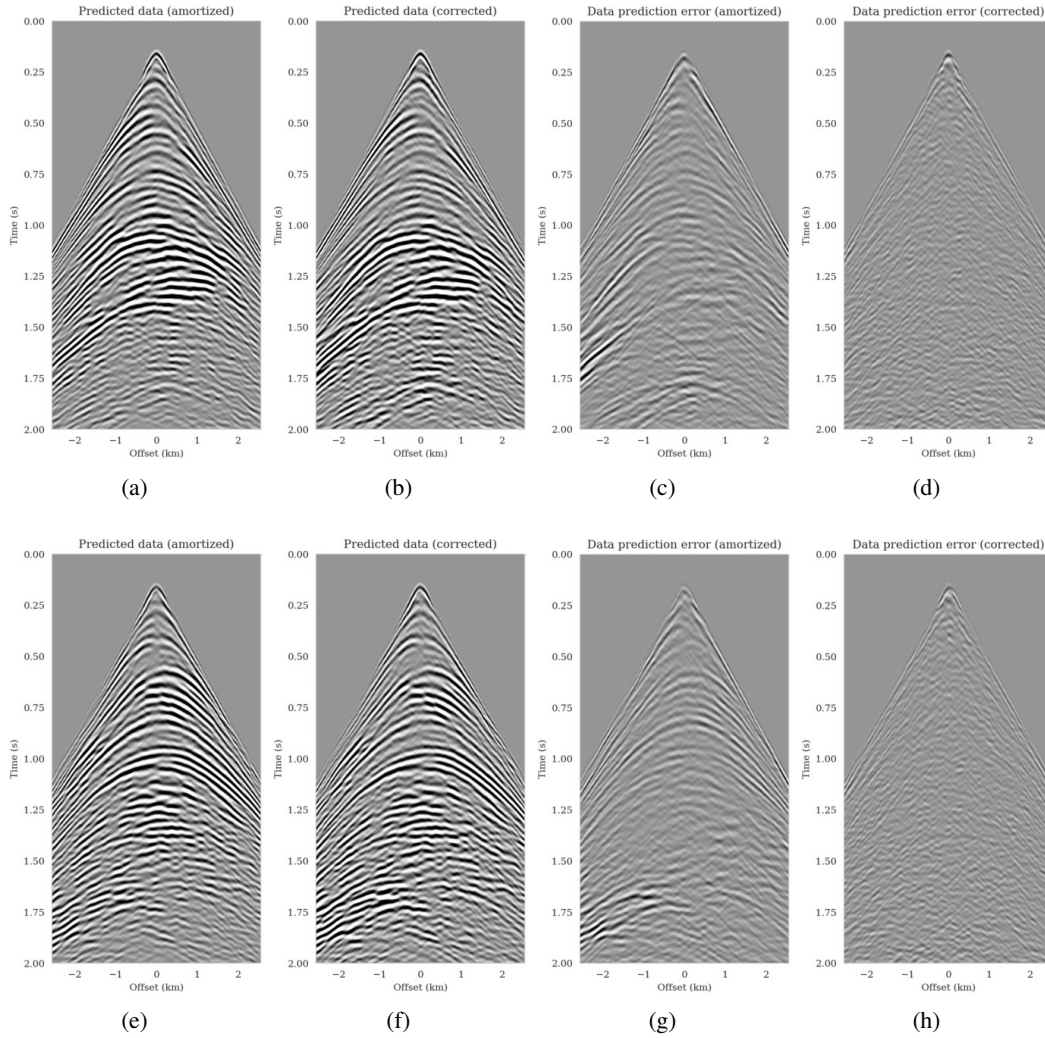


Figure 24: Quality control in data space. The first and second row correspond to experiments involving Figures 21a and 21b, respectively. (a) and (e) Predicted data before correction with SNRs 9.13 dB and 9.02 dB, respectively. (b) and (f) Predicted data after latent distribution correction with SNRs 15.27 dB and 14.52 dB, respectively. (c) and (g) Data residuals (before correction) associated with Figures 24a and 24e, respectively. (d) and (h) Data residuals (after correction) associated with Figures 24b and 24f, respectively.

account when dealing with large-scale nonlinear inverse problems. First, the parametric surrogate (conditional) distribution should be capable of approximating multi-modal densities. Invertible neural networks [41] are a suitable choice for this parameterization as they are known to be universal approximators [67, 68], meaning that invertible networks are capable of approximating a general class of diffeomorphisms. In other words, as long as there is a smooth invertible map between the latent space and the desired posterior distribution, invertible neural networks can be used to approximate the posterior distribution. The second consideration is regarding finding low-dimensional summary statistics [30] of observed data to avoid training the conditional normalizing flow over high dimensional data. In our linear seismic imaging example, we “summarized” seismic data (shot records) via the reverse-time migrated image. The conditions for which this summarization does not negatively bias the outcome of Bayesian inference (in the context of linear inverse problems) is described by Orozco et al. [81]. Further work is needed to design low-dimensional summary statistics for nonlinear inverse problems in order to successfully apply our framework to large-scale nonlinear inverse problems.

As far as seismic imaging is concerned, uncertainty can be attributed to two main sources [82–84]: (1) data errors, including measurement noise; (2) modeling errors, including linearization errors, which diminish with the accuracy of the background velocity model. The scope of our research is focused on the first source of uncertainty, and we will explore how variational inference models can be used to capture errors in background models in future work. In contrast to the problem highlighted in this paper, capturing the uncertainty caused by errors in the background model would require generating training data involving imaging experiments for a variety of plausible background velocity models, which would be computationally intensive. Recent developments in Fourier neural operators [85, 86] may prove to be useful in addressing this problem.

## 7 Conclusions

In high-dimensional inverse problems with computationally expensive forward modeling operators, Bayesian inference is challenging due to the cost of sampling the high-dimensional posterior distribution. The computational costs associated with the forward operator often limit the applicability of sampling the posterior distribution with traditional Markov chain Monte Carlo and variational inference methods. Added to the computational challenges are the difficulties associated with selecting a prior distribution that encodes our prior knowledge while not negatively biasing the outcome of Bayesian inference. Amortized variational inference addresses these challenges by incurring an offline initial training cost for a deep neural network that can approximate the posterior distribution for previously unseen observed data, distributed according to the training data distribution. When high-quality training data is readily available, and as long as there are no shifts in the data distribution, the pretrained network is capable of providing samples from the posterior distribution for previously unknown data virtually free of additional costs.

Unfortunately, in certain domains, such as geophysical inverse problems, where the structure of the Earth’s subsurface is unknown, it can be challenging to obtain a training dataset, e.g., a collection of images of the subsurface, which statistically captures the strong heterogeneity exhibited by the Earth’s subsurface. Furthermore, changes to the data generation process could negatively influence the quality of Bayesian inferences with amortized variational inferences due to generalization errors associated with neural networks. To address these challenges while exploiting the computational benefits of amortized variational inference, we proposed a data-specific, physics-based, and computationally cheap correction to the latent distribution of a conditional normalizing flow, pretrained to via an amortized variational inference objective. This correction involves solving a variational inference problem in the latent space of the pretrained conditional normalizing flow where we obtain a diagonal correction to the latent distribution such that the predicted posterior distribution more closely matches the desired posterior distribution.

Using a seismic imaging example, we demonstrate that the proposed latent distribution correction, at a cost of five reverse-time migrations, can be used to mitigate the effects of data distribution shifts, which includes changes in the forward model as well as the prior distribution. Our evaluation indicated improvements in seismic image quality, comparable to least squares imaging, after the latent distribution correction step, as well as estimate on the uncertainty of the image. We presented the pointwise standard deviation as a measure of uncertainty in the image, which indicated an increase in variability in complex geological areas and poorly illuminated areas. This approach will enable

uncertainty quantification in large-scale inverse problems, which otherwise would be computationally expensive to achieve.

## 8 Related material

The latent distribution correction optimization problem (equation 19) involves computing gradients of the composition of the forward operator and the pretrained conditional normalizing flow with respect to the latent variable. Computing this gradient requires actions of the forward operator and the adjoint of its Jacobian. In our numerical experiments, these operations involved solving wave equations. For maximal numerical performance, we use **JUDI** [87] to construct wave-equation solvers, which utilizes the just-in-time **Devito** [88, 89] compiler for the wave-equation based simulations. The invertible network architectures are implemented using **InvertibleNetworks.jl** [90], a memory-efficient framework for training invertible nets in Julia programming language. For more details on our implementation, please refer to our code on [GitHub](#).

## 9 Acknowledgments

This research was carried out with the support of Georgia Research Alliance and partners of the ML4Seismic Center.

## 10 Appendix A

### 10.1 Derivation of latent distribution correction objective function

The optimization problem for correcting the latent distribution involves minimizing the KL divergence between the Gaussian relaxation of the latent distribution  $N(\mathbf{z} \mid \boldsymbol{\mu}, \text{diag}(\mathbf{s})^2)$  and the shifted latent distribution  $p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}})$ , which is conditioned on an instance of out-of-distribution data  $\mathbf{y}_{\text{obs}} \sim \hat{p}_{\text{data}}(\mathbf{y})$ . This is an instance of non-amortized variational inference (see equations 4 and 5) defined over the latent variable:

$$\begin{aligned} & \arg \min_{\boldsymbol{\mu}, \mathbf{s}} \mathbb{KL} \left( N(\mathbf{z} \mid \boldsymbol{\mu}, \text{diag}(\mathbf{s})^2) \parallel p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}}) \right) \\ & = \arg \min_{\boldsymbol{\mu}, \mathbf{s}} \mathbb{E}_{\mathbf{z} \sim N(\mathbf{z} \mid \boldsymbol{\mu}, \text{diag}(\mathbf{s})^2)} \left[ -\log p_\phi(\mathbf{z} \mid \mathbf{y}_{\text{obs}}) + \log N(\mathbf{z} \mid \boldsymbol{\mu}, \text{diag}(\mathbf{s})^2) \right]. \end{aligned} \quad (21)$$

The above expression can be further simplified by rewriting the expectation with respect to  $N(\mathbf{z} \mid \boldsymbol{\mu}, \text{diag}(\mathbf{s})^2)$  as the expectation with respect to a standard Gaussian distribution, followed by an elementwise scaling by  $\mathbf{s}$  and a shift by  $\boldsymbol{\mu}$  [see reparameterization trick in 64], i.e.,

$$\arg \min_{\boldsymbol{\mu}, \mathbf{s}} \mathbb{E}_{\mathbf{z} \sim N(\mathbf{z} \mid \mathbf{0}, \mathbf{I})} \left[ -\log p_\phi(\mathbf{s} \odot \mathbf{z} + \boldsymbol{\mu} \mid \mathbf{y}_{\text{obs}}) + \log N(\mathbf{s} \odot \mathbf{z} + \boldsymbol{\mu} \mid \boldsymbol{\mu}, \text{diag}(\mathbf{s})^2) \right]. \quad (22)$$

The last term in the expectation in equation 22 is the log-density of a Gaussian distribution, which is equal to:

$$\begin{aligned} & \log N(\mathbf{s} \odot \mathbf{z} + \boldsymbol{\mu} \mid \boldsymbol{\mu}, \text{diag}(\mathbf{s})^2) \\ & = -\frac{D}{2} \log(2\pi) - \frac{1}{2} \log \left| \det \text{diag}(\mathbf{s})^2 \right| - \frac{1}{2} (\mathbf{s} \odot \mathbf{z} + \boldsymbol{\mu} - \boldsymbol{\mu})^\top \text{diag}(\mathbf{s})^{-2} (\mathbf{s} \odot \mathbf{z} + \boldsymbol{\mu} - \boldsymbol{\mu}) \\ & = -\log \left| \det \text{diag}(\mathbf{s}) \right| - \frac{1}{2} \|\mathbf{z}\|_2^2 + \text{const} \\ & = -\log \left| \det \text{diag}(\mathbf{s}) \right| + \text{const}. \end{aligned} \quad (23)$$

In the above equation,  $D$  is the dimension of  $\mathbf{z}$ , and the constants represent terms that are not a function of  $\boldsymbol{\mu}$  or  $\mathbf{s}$ . By inserting equation 23 into equation 22 we arrive at the following objective function for the latent distribution correction:

$$\arg \min_{\boldsymbol{\mu}, \mathbf{s}} \mathbb{E}_{\mathbf{z} \sim N(\mathbf{z} \mid \mathbf{0}, \mathbf{I})} \left[ -\log p_\phi(\mathbf{s} \odot \mathbf{z} + \boldsymbol{\mu} \mid \mathbf{y}_{\text{obs}}) - \log \left| \det \text{diag}(\mathbf{s}) \right| \right] \quad (24)$$

Finally, we use Bayes' rule and inserting the shifted latent density function from equation 17 to arrive at the objective function for latent distribution correction (equation 19):

$$\arg \min_{\boldsymbol{\mu}, \mathbf{s}} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})} \left[ \frac{1}{2\sigma^2} \sum_{i=1}^N \|\mathbf{y}_{\text{obs},i} - \mathcal{F}_i \circ f_{\phi}(\mathbf{s} \odot \mathbf{z} + \boldsymbol{\mu}; \mathbf{y}_{\text{obs}})\|_2^2 + \frac{1}{2} \|\mathbf{s} \odot \mathbf{z} + \boldsymbol{\mu}\|_2^2 - \log |\det \text{diag}(\mathbf{s})| \right]. \quad (25)$$

## References

- [1] Richard C Aster, Brian Borchers, and Clifford H Thurber. *Parameter Estimation and Inverse Problems*. Elsevier, 2018. doi: 10.1016/C2015-0-02458-3.
- [2] Albert Tarantola. *Inverse problem theory and methods for model parameter estimation*. SIAM, 2005. ISBN 978-0-89871-572-9. doi: 10.1137/1.9780898717921.
- [3] C. Robert and G. Casella. *Monte Carlo statistical methods*. Springer-Verlag, 2004.
- [4] James Martin, Lucas C. Wilcox, Carsten Burstedde, and OMAR Ghattas. A Stochastic Newton MCMC Method for Large-scale Statistical Inverse Problems with Application to Seismic Inversion. *SIAM Journal on Scientific Computing*, 34(3):A1460–A1487, 2012. URL <http://epubs.siam.org/doi/abs/10.1137/110845598>.
- [5] Zhilong Fang, Curt Da Silva, Rachel Kuske, and Felix J. Herrmann. Uncertainty quantification for inverse problems with weak partial-differential-equation constraints. *GEOPHYSICS*, 83(6): R629–R647, 2018. doi: 10.1190/geo2017-0824.1.
- [6] Ali Siahkoohi, Gabrio Rizzuti, and Felix J. Herrmann. Deep Bayesian inference for seismic imaging with tasks. *Geophysics*, 87(5), 6 2022. doi: 10.1190/geo2021-0666.1. URL <https://arxiv.org/abs/2110.04825>.
- [7] Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin. *Bayesian Data Analysis*. CRC press, 2013. doi: 10.1201/9780429258480. URL <http://www.stat.columbia.edu/~gelman/book/BDA3.pdf>.
- [8] Andrew Curtis and Anthony Lomax. Prior information, sampling distributions, and the curse of dimensionality. *Geophysics*, 66(2):372–378, 2001.
- [9] Max Welling and Yee Whye Teh. Bayesian Learning via Stochastic Gradient Langevin Dynamics. In *Proceedings of the 28th International Conference on Machine Learning, ICML'11*, pages 681–688, Madison, WI, USA, 2011. Omnipress. ISBN 9781450306195. doi: 10.5555/3104482.3104568. URL <https://dl.acm.org/doi/abs/10.5555/3104482.3104568>.
- [10] Felix J. Herrmann, Ali Siahkoohi, and Gabrio Rizzuti. Learned imaging with constraints and uncertainty quantification. In *Neural Information Processing Systems (NeurIPS) 2019 Deep Inverse Workshop*, 12 2019. URL <https://arxiv.org/pdf/1909.06473.pdf>.
- [11] Zeyu Zhao and Mrinal K Sen. A gradient based MCMC method for FWI and uncertainty analysis. In *89th Annual International Meeting, SEG*, pages 1465–1469. Expanded Abstracts, 2019. doi: 10.1190/segam2019-3216560.1.
- [12] M Kotsi, A Malcolm, and G Ely. Uncertainty quantification in time-lapse seismic imaging: a full-waveform approach. *Geophysical Journal International*, 222(2):1245–1263, 05 2020. ISSN 0956-540X. doi: 10.1093/gji/ggaa245.
- [13] Ali Siahkoohi, Gabrio Rizzuti, and Felix J. Herrmann. A deep-learning based bayesian approach to seismic imaging and uncertainty quantification. In *82nd EAGE Conference and Exhibition*. Extended Abstracts, 2020. doi: 10.3997/2214-4609.202010770.
- [14] Ali Siahkoohi, Gabrio Rizzuti, and Felix J. Herrmann. Uncertainty quantification in imaging and automatic horizon tracking—a Bayesian deep-prior based approach. In *90th Annual International Meeting, SEG*, pages 1636–1640. Expanded Abstracts, 9 2020. doi: 10.1190/segam2020-3417560.1.

- [15] Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An Introduction to Variational Methods for Graphical Models. *Machine Learning*, 37(2):183–233, 1999. doi: 10.1023/A:1007665907178.
- [16] Martin J Wainwright and Michael Irwin Jordan. *Graphical models, exponential families, and variational inference*. Now Publishers Inc, 2008.
- [17] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. volume 37 of *Proceedings of Machine Learning Research*, pages 1530–1538. PMLR, 07–09 Jul 2015. URL <http://proceedings.mlr.press/v37/rezende15.html>.
- [18] Qiang Liu and Dilin Wang. Stein Variational Gradient Descent: A General Purpose Bayesian Inference Algorithm. In *Advances in Neural Information Processing Systems*, volume 29, pages 2378–2386. Curran Associates, Inc., 2016. URL <https://proceedings.neurips.cc/paper/2016/file/b3ba8f1bee1238a2f37603d90b58898d-Paper.pdf>.
- [19] Gabrio Rizzuti, Ali Siahkoobi, Philipp A. Witte, and Felix J. Herrmann. Parameterizing uncertainty by deep invertible networks, an application to reservoir characterization. In *90th Annual International Meeting, SEG*, pages 1541–1545, 09 2020. doi: 10.1190/segam2020-3428150.1.
- [20] Xin Zhang and Andrew Curtis. Seismic tomography using variational inference methods. *Journal of Geophysical Research: Solid Earth*, 125(4):e2019JB018589, 2020.
- [21] Malte Tölle, Max-Heinrich Laves, and Alexander Schlaefer. A Mean-Field Variational Inference Approach to Deep Image Prior for Inverse Problems in Medical Imaging. In *Medical Imaging with Deep Learning*, 2021.
- [22] Dongzhuo Li, Huseyin Denli, Cody MacDonald, Kyle Basler-Reeder, Anatoly Baumstein, and Jacquelyn Daves. Multiparameter geophysical reservoir characterization augmented by generative networks. In *First International Meeting for Applied Geoscience & Energy*, pages 1364–1368. Society of Exploration Geophysicists, 2021.
- [23] Dongzhuo Li. Differentiable gaussianization layers for inverse problems regularized by deep generative models. *arXiv preprint arXiv:2112.03860*, 2022.
- [24] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- [25] Xin Zhang, Muhammad Atif Nawaz, Xuebin Zhao, and Andrew Curtis. An introduction to variational inference in geophysical inverse problems. In *Inversion of Geophysical Data*, pages 73–140. Elsevier, 2021. doi: 10.1016/bs.agph.2021.06.003. URL <https://doi.org/10.1016%2Fbs.agph.2021.06.003>.
- [26] Yoon Kim, Sam Wiseman, Andrew Miller, David Sontag, and Alexander Rush. Semi-amortized variational autoencoders. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2678–2687. PMLR, 10–15 Jul 2018. URL <http://proceedings.mlr.press/v80/kim18e.html>.
- [27] Ricardo Baptista, Olivier Zahm, and Youssef Marzouk. An adaptive transport framework for joint and conditional density estimation. *arXiv preprint arXiv:2009.10303*, 2020. URL <https://arxiv.org/abs/2009.10303>.
- [28] Jakob Kruse, Gianluca Detommaso, Robert Scheichl, and Ullrich Köthe. HINT: Hierarchical Invertible Neural Transport for Density Estimation and Bayesian Inference. *Proceedings of AAAI-2021*, 2021. URL <https://arxiv.org/pdf/1905.10687.pdf>.
- [29] Nikola Kovachki, Ricardo Baptista, Bamdad Hosseini, and Youssef Marzouk. Conditional Sampling With Monotone GANs, 2021.
- [30] Stefan T. Radev, Ulf K. Mertens, Andreas Voss, Lynton Ardizzone, and Ullrich Köthe. BayesFlow: Learning complex stochastic models with invertible neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 33(4):1452–1466, 2022. doi: 10.1109/TNNLS.2020.3042395.

- [31] Ali Siahkoohi, Gabrio Rizzuti, Mathias Louboutin, Philipp Witte, and Felix J. Herrmann. Preconditioned training of normalizing flows for variational inference in inverse problems. In *3rd Symposium on Advances in Approximate Bayesian Inference*, 1 2021. URL <https://openreview.net/pdf?id=P9m1sMaNQ8T>.
- [32] Yuxiao Ren, Philipp A. Witte, Ali Siahkoohi, Mathias Louboutin, and Felix J. Herrmann. Seismic Velocity Inversion and Uncertainty Quantification Using Conditional Normalizing Flows. In *American Geophysical Union (AGU) Fall Meeting*, 12 2021. URL <https://agu.confex.com/agu/fm21/meetingapp.cgi/Paper/815883>.
- [33] Ali Siahkoohi and Felix J Herrmann. Learning by example: fast reliability-aware seismic imaging with normalizing flows. In *First International Meeting for Applied Geoscience & Energy*, pages 1580–1585. Expanded Abstracts, 2021. doi: 10.1190/segam2021-3581836.1.
- [34] AmirEhsan Khorashadizadeh, Konik Kothari, Leonardo Salsi, Ali Aghababaei Harandi, Maarten de Hoop, and Ivan Dokmani'c. Conditional Injective Flows for Bayesian Imaging. *arXiv preprint arXiv:2204.07664*, 2022.
- [35] Rafael Orozco, Ali Siahkoohi, Gabrio Rizzuti, Tristan van Leeuwen, and Felix Johan Herrmann. Photoacoustic imaging with conditional priors from normalizing flows. In *NeurIPS 2021 Workshop on Deep Learning and Inverse Problems*, 2021. URL <https://openreview.net/forum?id=woi1OTvROO1>.
- [36] Ali Siahkoohi, Rafael Orozco, Gabrio Rizzuti, and Felix J. Herrmann. Wave-equation based inversion with amortized variational bayesian inference. In *EAGE Deep learning for seismic processing: Investigating the foundations workshop*, 6 2022. URL <https://arxiv.org/abs/2203.15881>.
- [37] Amirhossein Taghvaei and Bamdad Hosseini. An optimal transport formulation of bayes' law for nonlinear filtering algorithms. *arXiv preprint arXiv:2203.11869*, 2022.
- [38] Jian Sun, Kristopher A Innanen, and Chao Huang. Physics-guided deep learning for seismic inversion with hybrid training and uncertainty analysis. *Geophysics*, 86(3):R303–R317, 2021.
- [39] Peng Jin, Xitong Zhang, Yinpeng Chen, Sharon X Huang, Zicheng Liu, and Youzuo Lin. Unsupervised learning of full-waveform inversion: Connecting CNN and partial differential equation in a loop. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=izvbwgBic9q>.
- [40] Marvin Schmitt, Paul-Christian Bürkner, Ullrich Köthe, and Stefan T Radev. Detecting model misspecification in amortized bayesian inference with neural networks. *arXiv preprint arXiv:2112.08866*, 2021.
- [41] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using Real NVP. In *International Conference on Learning Representations, ICLR*, 2016. URL <http://arxiv.org/abs/1605.08803>.
- [42] Muhammad Asim, Max Daniels, Oscar Leong, Ali Ahmed, and Paul Hand. Invertible generative models for inverse problems: mitigating representation error and dataset bias. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 399–409. PMLR, 07 2020. URL <http://proceedings.mlr.press/v119/asim20a.html>.
- [43] Matthew D Parno and Youssef M Marzouk. Transport Map Accelerated Markov Chain Monte Carlo. *SIAM/ASA Journal on Uncertainty Quantification*, 6(2):645–682, 2018. doi: 10.1137/17M1134640.
- [44] Benjamin Peherstorfer and Youssef Marzouk. A transport-based multifidelity preconditioner for Markov chain Monte Carlo. *Advances in Computational Mathematics*, 45(5-6):2321–2348, 2019.

- [45] Ashish Bora, Ajil Jalal, Eric Price, and Alexandros G. Dimakis. Compressed sensing using generative models. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 537–546. PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/bora17a.html>.
- [46] Anna Andriele, Nando Farchmin, Paul Hagemann, Sebastian Heidenreich, Victor Soltwisch, and Gabriele Steidl. Invertible neural networks versus mcmc for posterior reconstruction in grazing incidence x-ray fluorescence. In *International Conference on Scale Space and Variational Methods in Computer Vision*, pages 528–539. Springer, 2021.
- [47] Xuebin Zhao, Andrew Curtis, and Xin Zhang. Bayesian seismic tomography using normalizing flows. *Geophysical Journal International*, 228(1):213–239, 07 2021. ISSN 0956-540X. doi: 10.1093/gji/ggab298. URL <https://doi.org/10.1093/gji/ggab298>.
- [48] Xin Zhang and Andrew Curtis. Bayesian geophysical inversion using invertible neural networks. *Journal of Geophysical Research: Solid Earth*, 126(7):e2021JB022320, 2021.
- [49] Xuebin Zhao, Andrew Curtis, and Xin Zhang. Interrogating subsurface structures using probabilistic tomography: an example assessing the volume of irish sea basins. *Journal of Geophysical Research: Solid Earth*, 127(4):e2022JB024098, 2022.
- [50] Konik Kothari, AmirEhsan Khorashadizadeh, Maarten de Hoop, and Ivan Dokmanic. Trumpets: Injective Flows for Inference and Inverse Problems, 2021.
- [51] Jonas Adler and Ozan Öktem. Deep Bayesian Inversion. *arXiv preprint arXiv:1811.05910*, 2018. URL <https://arxiv.org/abs/1811.05910>.
- [52] Jay Whang, Erik Lindgren, and Alex Dimakis. Composing normalizing flows for inverse problems. In *International Conference on Machine Learning*, pages 11158–11169. PMLR, 2021.
- [53] Alberto Malinverno and Victoria A Briggs. Expanded uncertainty quantification in inverse problems: Hierarchical bayes and empirical bayes. *GEOPHYSICS*, 69(4):1005–1016, 2004. doi: 10.1190/1.1778243.
- [54] Alberto Malinverno and Robert L Parker. Two ways to quantify uncertainty in geophysical inverse problems. *GEOPHYSICS*, 71(3):W15–W27, 2006. doi: 10.1190/1.2194516.
- [55] Anandaroop Ray, Sam Kaplan, John Washbourne, and Uwe Albertin. Low frequency full waveform seismic inversion within a tree based Bayesian framework. *Geophysical Journal International*, 212(1):522–542, 10 2017. ISSN 0956-540X. doi: 10.1093/gji/ggx428. URL <https://doi.org/10.1093/gji/ggx428>.
- [56] Georgia K Stuart, Susan E Minkoff, and Felipe Pereira. A two-stage Markov chain Monte Carlo method for seismic inversion and uncertainty quantification. *GEOPHYSICS*, 84(6):R1003–R1020, 11 2019. doi: 10.1190/geo2018-0893.1.
- [57] Chunyuan Li, Changyou Chen, David Carlson, and Lawrence Carin. Preconditioned Stochastic Gradient Langevin Dynamics for Deep Neural Networks. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, AAAI’16, pages 1788–1794. AAAI Press, 2016. doi: 10.5555/3016100.3016149. URL <https://dl.acm.org/doi/abs/10.5555/3016100.3016149>.
- [58] Christopher M Bishop. *Pattern Recognition and Machine Learning*. Springer-Verlag New York, 2006.
- [59] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- [60] Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4): 1574–1609, 2009.



- [61] Tijmen Tieleman and Geoffrey Hinton. Lecture 6.5-RMSprop: Divide the gradient by a running average of its recent magnitude. [https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture\\_slides\\_lec6.pdf](https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf), 2012.
- [62] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. URL <https://arxiv.org/pdf/1412.6980.pdf>.
- [63] Cédric Villani. *Optimal transport: old and new*. Springer-Verlag Berlin Heidelberg, 2009. doi: 10.1007/978-3-540-71050-9.
- [64] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014.
- [65] Kyle Cranmer, Johann Brehmer, and Gilles Louppe. The frontier of simulation-based inference. *Proceedings of the National Academy of Sciences*, 117(48):30055–30062, 2020.
- [66] Alexander Lavin, Hector Zenil, Brooks Paige, David Krakauer, Justin Gottschlich, Tim Mattson, Anima Anandkumar, Sanjay Choudry, Kamil Rocki, Atılım Güneş Baydin, et al. Simulation intelligence: Towards a new generation of scientific methods. *arXiv preprint arXiv:2112.03235*, 2021.
- [67] Takeshi Teshima, Isao Ishikawa, Koichi Tojo, Kenta Oono, Masahiro Ikeda, and Masashi Sugiyama. Coupling-based invertible neural networks are universal diffeomorphism approximators. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 3362–3373. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/2290a7385ed77cc5592dc2153229f082-Paper.pdf>.
- [68] Isao Ishikawa, Takeshi Teshima, Koichi Tojo, Kenta Oono, Masahiro Ikeda, and Masashi Sugiyama. Universal approximation property of invertible neural networks. *arXiv preprint arXiv:2204.07415*, 2022. URL <https://arxiv.org/abs/2204.07415>.
- [69] J. E. Gubernatis, E. Domany, J. A. Krumhansl, and M. Huberman. The Born approximation in the theory of the scattering of elastic waves by flaws. *Journal of Applied Physics*, 48(7): 2812–2819, 1977. doi: 10.1063/1.324142.
- [70] Gilles Lambaré, Jean Virieux, Raul Madariaga, and Side Jin. Iterative asymptotic inversion in the acoustic approximation. *Geophysics*, 57(9):1138–1154, 1992.
- [71] Gerard T Schuster. Least-squares cross-well migration. In *SEG Technical Program Expanded Abstracts 1993*, pages 110–113. Society of Exploration Geophysicists, 1993.
- [72] Tamas Nemeth, Chengjun Wu, and Gerard T Schuster. Least-squares migration of incomplete reflection data. *GEOPHYSICS*, 64(1):208–221, 1999. doi: 10.1190/1.1444517.
- [73] Veritas. Parihaka 3D Marine Seismic Survey - Acquisition and Processing Report. Technical Report New Zealand Petroleum Report 3460, New Zealand Petroleum & Minerals, Wellington, 2005.
- [74] WesternGeco. Parihaka 3D PSTM Final Processing Report. Technical Report New Zealand Petroleum Report 4582, New Zealand Petroleum & Minerals, Wellington, 2012.
- [75] Jonas Adler, Sebastian Lunz, Olivier Verdier, Carola-Bibiane Schönlieb, and Ozan Öktem. Task adapted reconstruction for inverse problems. *Inverse Problems*, 38(7):075006, may 2022. doi: 10.1088/1361-6420/ac28ec. URL <https://doi.org/10.1088/1361-6420/ac28ec>.
- [76] Ruitao Liu, Arijit Chakrabarti, Tapas Samanta, Jayanta K Ghosh, and Malay Ghosh. On divergence measures leading to jeffreys and other reference priors. *Bayesian Analysis*, 9(2): 331–370, 2014.
- [77] Yanchao Yang and Stefano Soatto. Conditional Prior Networks for Optical Flow. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 271–287, 2018.

- [78] Chen Zeno, Itay Golan, Elad Hoffer, and Daniel Soudry. Task Agnostic Continual Learning Using Online Variational Bayes. *arXiv preprint arXiv:1803.10123*, 2018.
- [79] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In *Proceedings of the 27th International Conference on Neural Information Processing Systems, NIPS'14*, pages 3320–3328, 2014. URL <http://dl.acm.org/citation.cfm?id=2969033.2969197>.
- [80] Xin Zhang, Angus Lomas, Muhong Zhou, York Zheng, and Andrew Curtis. 3d bayesian variational full waveform inversion, 2022. URL <https://arxiv.org/abs/2210.03613>.
- [81] Rafael Orozco, Ali Siahkoohi, Gabrio Rizzuti, Tristan van Leeuwen, and Felix J. Herrmann. Adjoint operators enable fast and amortized machine learning based bayesian uncertainty quantification. In *SPIE Medical Imaging Conference*, 02 2023. URL [https://slim.gatech.edu/Publications/Public/Conferences/SPIE/2023/orozco2023SPIEadjoint/SPIE\\_2022\\_adjoint.html](https://slim.gatech.edu/Publications/Public/Conferences/SPIE/2023/orozco2023SPIEadjoint/SPIE_2022_adjoint.html). (SPIE, San Diego).
- [82] Pierre Thore, Arben Shtuka, Magali Lecour, Taoufik Ait-Ettajer, and Richard Cognot. Structural uncertainties: Determination, management, and applications. *Geophysics*, 67(3):840–852, 2002.
- [83] Konstantin Osypov, Yi Yang, Aimé Fournier, Natalia Ivanova, Ran Bachrach, Can Evren Yarman, Yu You, Dave Nichols, and Marta Woodward. Model-uncertainty quantification in seismic tomography: method and applications. *Geophysical Prospecting*, 61(6-Challenges of Seismic Imaging and Inversion Devoted to Goldin):1114–1134, 2013.
- [84] Gregory Ely, Alison Malcolm, and Oleg V. Poliannikov. Assessing uncertainties in velocity models and images with a fast nonlinear uncertainty quantification method. *GEOPHYSICS*, 83(2):R63–R75, 2018. doi: 10.1190/geo2017-0321.1. URL <https://doi.org/10.1190/geo2017-0321.1>.
- [85] Zongyi Li, Nikola Borislavov Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Kaushik Bhattacharya, Andrew M. Stuart, and Anima Anandkumar. Fourier Neural Operator for Parametric Partial Differential Equations. In *9th International Conference on Learning Representations*. OpenReview.net, 2021. URL <https://openreview.net/forum?id=c8P9NQVtmnO>.
- [86] Ali Siahkoohi, Mathias Louboutin, and Felix J. Herrmann. Velocity continuation with Fourier neural operators for accelerated uncertainty quantification. In *2nd International Meeting for Applied Geoscience & Energy*, 2022.
- [87] Philipp A. Witte, Mathias Louboutin, Navjot Kukreja, Fabio Luporini, Michael Lange, Gerard J. Gorman, and Felix J. Herrmann. A large-scale framework for symbolic implementations of seismic inversion algorithms in julia. *GEOPHYSICS*, 84(3):F57–F71, 2019. doi: 10.1190/geo2018-0174.1. URL <https://doi.org/10.1190/geo2018-0174.1>.
- [88] F. Luporini, M. Lange, M. Louboutin, N. Kukreja, J. Hükelheim, C. Yount, P. Witte, P. H. J. Kelly, F. J. Herrmann, and G. J. Gorman. Architecture and performance of devito, a system for automated stencil computation. *CoRR*, abs/1807.03032, jul 2018. URL <http://arxiv.org/abs/1807.03032>.
- [89] M. Louboutin, M. Lange, F. Luporini, N. Kukreja, P. A. Witte, F. J. Herrmann, P. Velesko, and G. J. Gorman. Devito (v3.1.0): an embedded domain-specific language for finite differences and geophysical exploration. *Geoscientific Model Development*, 12(3):1165–1187, 2019. doi: 10.5194/gmd-12-1165-2019. URL <https://www.geosci-model-dev.net/12/1165/2019/>.
- [90] Philipp Witte, Gabrio Rizzuti, Mathias Louboutin, Ali Siahkoohi, Felix Herrmann, and Bas Peters. InvertibleNetworks.jl: A Julia framework for invertible neural networks, March 2021. URL <https://doi.org/10.5281/zenodo.4610118>.