

# A-optimal encoding weights for nonlinear inverse problems, with application to the Helmholtz inverse problem

Benjamin Crestel<sup>1</sup>, Alen Alexanderian<sup>2</sup>, Georg Stadler<sup>3</sup>  
and Omar Ghattas<sup>1,4</sup>

<sup>1</sup>Institute for Computational Engineering & Sciences, The University of Texas at Austin, Austin, TX, USA

<sup>2</sup>Department of Mathematics, North Carolina State University, Raleigh, NC, USA

<sup>3</sup>Courant Institute of Mathematical Sciences, New York University, New York, NY, USA

<sup>4</sup>Department of Geological Sciences and Department of Mechanical Engineering, The University of Texas at Austin, Austin, TX, USA

E-mail: [crestel@ices.utexas.edu](mailto:crestel@ices.utexas.edu), [alexanderian@ncsu.edu](mailto:alexanderian@ncsu.edu),  
[stadler@cims.nyu.edu](mailto:stadler@cims.nyu.edu) and [omar@ices.utexas.edu](mailto:omar@ices.utexas.edu)

January 2016

**Abstract.** The computational cost of solving an inverse problem governed by PDEs, using multiple experiments, increases linearly with the number of experiments. A recently proposed method to decrease this cost uses only a small number of random linear combinations of all experiments for solving the inverse problem. This approach applies to inverse problems where the PDE solution depends linearly on the right-hand side function that models the experiment. As this method is stochastic in essence, the quality of the obtained reconstructions can vary, in particular when only a small number of combinations are used. We develop a Bayesian formulation for the definition and computation of encoding weights that lead to a parameter reconstruction with the least uncertainty. We call these weights A-optimal encoding weights. Our framework applies to inverse problems where the governing PDE is nonlinear with respect to the inversion parameter field. We formulate the problem in infinite dimensions and follow the optimize-then-discretize approach, devoting special attention to the discretization and the choice of numerical methods in order to achieve a computational cost that is independent of the parameter discretization. We elaborate our method for a Helmholtz inverse problem, and derive the adjoint-based expressions for the gradient of the objective function of the optimization problem for finding the A-optimal encoding weights. The proposed method is potentially attractive for real-time monitoring applications, where one can invest the effort to compute optimal weights offline, to later solve an inverse problem repeatedly, over time, at a fraction of the initial cost.

*Keywords:* source encoding, Bayesian nonlinear inverse problem, A-optimal experimental design, randomized trace estimator, Helmholtz equation.

Submitted to: *Inverse Problems*

## 1. Introduction

Inverse problems are ubiquitous in science and engineering. They arise whenever one attempts to infer parameters  $m$  from indirect observations  $\mathbf{d}$  and from a mathematical model—the parameter-to-observable map,  $\mathcal{F}(\cdot)$ —for the physical phenomenon that relates  $m$  and  $\mathbf{d}$ . When available, it is common to use observations obtained from different experiments to improve the quality of the parameter estimation. Suppose  $N_s$  experiments are conducted, indexed by  $i \in \{1, \dots, N_s\}$ . The  $i$ -th experiment results in observations  $\mathbf{d}_i$  and the corresponding parameter-to-observable map is denoted by  $\mathcal{F}_i(m)$ . Following a deterministic approach to this inverse problem results in the nonlinear least-squares minimization problem

$$\min_m \left\{ \frac{1}{2N_s} \sum_{i=1}^{N_s} \|\mathcal{F}_i(m) - \mathbf{d}_i\|^2 + \mathcal{R}(m) \right\}, \quad (1)$$

where  $\mathcal{R}$  is an appropriate regularization operator to cope with the ill-posedness that is common for many inverse problems.

Nonlinear optimization problems such as (1) can only be solved iteratively, which requires the availability of first (and ideally, also second) derivatives of the functional in (1) with respect to  $m$ . For an important class of inverse problems, the parameter-to-observable map involves the solution of a partial differential equation (PDE). This means that the evaluation of  $\mathcal{F}_i(m)$  entails the solution  $u_i$  of a PDE, and this  $u_i$  is usually restricted by an observation operator  $B$  to a subset of the domain (e.g., points), where observations are available. In this work, we make the assumption that the different experiments correspond to different right-hand sides  $f_i$  of this PDE. Moreover, this PDE must be linear with respect to the solution  $u_i$ , and both the PDE operator as well as the observation operator  $B$  must be the same for all experiments.

When the  $i$ -th experiment corresponds to a forcing term  $f_i$ , the parameter-to-observable map is given by  $\mathcal{F}_i(m) = Bu_i$ , where  $\mathcal{A}(m)u_i = f_i$  with  $\mathcal{A}(m)$  denoting the linear PDE-operator that may depend nonlinearly on  $m$ . Note that the governing PDE can be stationary or time-dependent. Adjoint methods allow to compute derivatives of the objective in (1) efficiently [1]. For instance, the computation of the gradient of the objective in (1) requires solving  $N_s$  forward and associated adjoint PDEs. Similar computational costs are associated with the application of the Hessian operator to vectors, such that the overall computational cost of solving (1), which is dominated by PDE solves with the operator  $\mathcal{A}(m)$ , grows (at least) linearly with the number of experiments  $N_s$ . In some important inverse problems,  $N_s$  is large (e.g., several thousand), such that these computations are expensive or even infeasible.

There have been some recent breakthroughs to address this computational bottleneck using the concept of random source encoding, sometimes also referred to as simultaneous random sources [2, 3]. A mathematical justification of this approach is given in the seminal paper [4], and is summarized in section 2. In [5], the authors employed a similar idea to encode the observations in inverse problems with large amount of data. The main idea of random source encoding is to replace the data generated by each individual experiment with a small number,  $N_w \ll N_s$ , of linear combinations of the data; the weights of these linear combinations,  $\mathbf{w}^i = [w_1^i, \dots, w_{N_s}^i]^T$ , are called encoding weights. Due to our linearity assumptions, this linear combination of data corresponds to the same linear combination of experiments, i.e., we can define encoded parameter-to-observable maps  $\mathcal{F}(\mathbf{w}^i; m)$ ,  $i = 1, \dots, N_w$ , as

follows

$$\mathcal{F}(\mathbf{w}^i; m) := \sum_{j=1}^{N_s} w_j^i \mathcal{F}_j(m) = B \left( \sum_{j=1}^{N_s} w_j^i u_j \right). \quad (2)$$

Observe that  $\sum_{j=1}^{N_s} w_j^i u_j$  can be computed by solving the *single* PDE

$$\mathcal{A}(m) \left( \sum_{j=1}^{N_s} w_j^i u_j \right) = \left( \sum_{j=1}^{N_s} w_j^i f_j \right).$$

Replacing the individual experiments with encoded experiments results in an inverse problem with lower computational complexity. The hope is that these linear combinations still carry most of the information contained in the individual experiments. As mentioned above, the source encoding method hinges on the linearity of the PDE describing the underlying physical phenomenon, such that the observables depends linearly on the forcing term. Additionally, the unicity of the observation operator  $B$  is necessary, but this requirement can be weakened in certain situations, e.g., if data from some experiments is missing [6].

The method of random source encoding, stochastic in essence, suffers from a few limitations. The key idea of the random source encoding approach is the conversion of the deterministic optimization (1) into a stochastic optimization problem. The expectation to be minimized is then approximated using a Monte-Carlo technique (see [4] or section 2). To reduce the computational cost of solving the inverse problem, one would like to choose the number of samples used in this Monte-Carlo approximation small. A small number of samples translates into a large variance for the Monte-Carlo estimator of the expectation. In practice, this manifests itself in large differences in the reconstructions obtained with different samples of encoding weights. An approach to remedy that difficulty is to select the weights deterministically [7, 8]. In particular, in [7], the author considers to select the weights that generate the greatest improvement from the current reconstruction, but the results are inconclusive. In [8], the authors choose the weights that minimize the expected medium misfit in the case of a discrete linear inverse problem, which is related to the approach we follow in this paper.

*Contributions* The main contributions of this article are as follows: (1) Drawing from recent developments in optimal experimental design (OED) for high- or infinite-dimensional inverse problems [9, 10, 11, 12], and following a Bayesian view of inverse problems, we develop a method for the computation of encoding weights that lead to a parameter reconstruction with the least uncertainty—as measured by the average of the posterior variance. We refer to these (deterministic) weights as *A-optimal encoding weights*, a nomenclature motivated by the use of the A-optimal experimental design criterion from OED theory [13]. (2) The method we propose extends the work in [8] by addressing inverse problems with nonlinear parameter-to-observable maps, and allows for infinite-dimensional parameters. The infinite-dimensional formulation has two main advantages: (a) the use of weak forms facilitates the derivation of adjoint-based expressions for the gradient of the objective function to compute the A-optimal encoding weights; (b) it allows us to follow the optimize-then-discretize approach, which, along with devoting special attention to the discretization of the formulation and the choice of the numerical methods employed, helps control the

computational cost independently of the parameter discretization. (3) We elaborate our method for the Helmholtz inverse problem and derive the adjoint-based gradient of the optimization problem for finding the A-optimal encoding weights. We also analyze the computational cost—in terms of Helmholtz PDE solves—of objective and gradient evaluation for this optimization problem. For this Helmholtz problem, we present an extensive numerical study and discuss the potential and pitfalls of our approach.

*Paper overview* The rest of this article is organized as follows. In section 2, we provide an overview of the method of random source encoding. We also introduce notation that we will carry throughout the paper. In section 3, we summarize elements of Bayesian inverse problems and introduce approximations to the posterior covariance in function space. The framework for the A-optimal encoding weights is presented in section 4. In section 5, we elaborate our formulation for the Helmholtz inverse problem. We derive adjoint-based expressions for the gradient of the A-optimal objective function, and analyze computational cost of evaluating the objective function and its gradient. Numerical results are presented in section 6, and we provide some concluding remarks in section 7.

## 2. Random source encoding

In this section, we review the method of random source encoding, and introduce notation and terminology used throughout this article. We seek to infer a parameter field  $m \in \mathcal{V}$  where  $\mathcal{V}$  is an infinite-dimensional Hilbert space of functions defined over the domain  $\mathcal{D} \subset \mathbb{R}^d$  ( $d = 2, 3$ ); a typical choice is  $\mathcal{V} := L^2(\mathcal{D})$ . The parameter-to-observable map is denoted by  $\mathcal{F}_i : \mathcal{V} \rightarrow \mathbb{R}^q$ . Let us assume that  $u_i$  solves the PDE  $\mathcal{A}(m)u_i = f_i$  and that all experiments  $i = 1, \dots, N_s$  share a common observation operator  $B$ , where  $Bu_i \in \mathbb{R}^q$ . We then write each parameter-to-observable map as  $\mathcal{F}_i(m) = Bu_i$ . The right-hand side source  $f_i$  characterizes the  $i$ -th experiment. To apply source encoding, we require the parameter-to-observable map to be linear with respect to the source terms, which led us to introduce the encoded parameter-to-observable maps (2).

In [4] the authors give a mathematical justification of the idea of random source encoding for a discrete problem and we follow their argument, here, for an inverse problem formulated in function space. We gather all  $\mathcal{F}_i(m)$  (resp.  $\mathbf{d}_i$ ) in the columns of a matrix  $\mathbf{F}(m)$  (resp.  $\mathbf{D}^e$ ) and call the data misfit matrix  $\mathcal{S}(m) := \mathbf{F}(m) - \mathbf{D}^e$ . Ignoring the regularization term for now, the inverse problem can be written as,  $\min_{m \in \mathcal{V}} \left\{ \|\mathcal{S}(m)\|_F^2 \right\}$ , where  $\|\cdot\|_F$  is the Frobenius norm [14]. Note that  $\|\mathcal{S}(m)\|_F^2 = \text{trace}(\mathcal{S}(m)^T \mathcal{S}(m))$ , which can be approximated efficiently using randomized trace estimators [15, 16]. Indeed, for random vectors  $\mathbf{z}$  with mean zero and identity covariance matrix one finds that,  $\text{trace}(\mathcal{S}(m)^T \mathcal{S}(m)) = \mathbb{E}_{\mathbf{z}}(\|\mathcal{S}(m)\mathbf{z}\|_2^2)$ . Typical choices of distribution for  $\mathbf{z}$  include the Rademacher distribution, where samples take values  $\pm 1$  with probability  $1/2$ , and the standard normal distribution  $\mathcal{N}(0, \mathbf{I}_{N_s})$ . Among other possible choices we mention the discrete distribution that takes values  $\pm\sqrt{3}$  with probability  $1/6$  and 0 otherwise, or the uniform spherical distribution on a sphere of radius  $\sqrt{N_s}$  that we denote  $\mathcal{U}_s(\sqrt{N_s})$ ; the fact that  $\mathcal{U}_s(\sqrt{N_s})$  has identity covariance matrix can be shown using results from [17], along with the observation that  $\tilde{z} \sim \mathcal{U}_s(\sqrt{N_s})$  iff  $\tilde{z} = \sqrt{N_s}(\mathbf{z}/\|\mathbf{z}\|)$  with  $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I}_{N_s})$ . We now write the data-misfit term as an expectation, i.e.,  $\|\mathbf{F}(m) - \mathbf{D}^e\|_F^2 = \mathbb{E}_{\mathbf{z}}(\|(\mathbf{F}(m) - \mathbf{D}^e)\mathbf{z}\|^2)$ , leading

to the stochastic optimization problem

$$\min_{m \in \mathcal{V}} \left\{ \mathbb{E}_{\mathbf{z}} (\|(\mathbf{F}(m) - \mathbf{D}^e)\mathbf{z}\|^2) \right\}.$$

There exist two main techniques to solve these types of problems [18]. Using stochastic average approximation (SAA), one approximates the cost functional with a Monte-Carlo-type approach before solving a deterministic optimization problem, i.e., for fixed samples  $\mathbf{z}_i$  ones solves

$$\mathbb{E}_{\mathbf{z}} (\|(\mathbf{F}(m) - \mathbf{D}^e)\mathbf{z}\|^2) \approx \frac{1}{M} \sum_{i=1}^M \|(\mathbf{F}(m) - \mathbf{D}^e)\mathbf{z}_i\|^2.$$

In an alternative approach called stochastic approximation (SA), one re-samples the random vector  $\mathbf{z}$  at each step of the iteration.

We now specify the source-encoded equivalent of (1). Given  $N_w$  encoding weights  $\mathbf{w} = (\mathbf{w}^1, \dots, \mathbf{w}^{N_w})$ , where each  $\mathbf{w}^i \in \mathbb{R}^{N_s}$ , we define the encoded data  $\mathbf{d}(\mathbf{w}^i) := \sum_{j=1}^{N_s} w_j^i \mathbf{d}_j$ , the encoded right-hand side  $f(\mathbf{w}^i) := \sum_{j=1}^{N_s} w_j^i f_j$ , and encoded parameter-to-observable maps  $\mathcal{F}(\mathbf{w}^i; m) = \sum_{j=1}^{N_s} w_j^i \mathcal{F}_j(m)$ . The parameter field  $m_c(\mathbf{w})$  reconstructed using the  $N_w$  encoded sources is then defined as

$$m_c(\mathbf{w}) = \arg \min_{m \in \mathcal{V}} \left\{ \frac{1}{2N_w} \sum_{i=1}^{N_w} \|\mathcal{F}(\mathbf{w}^i; m) - \mathbf{d}(\mathbf{w}^i)\|^2 + \mathcal{R}(m) \right\}. \quad (3)$$

Due to the assumptions on  $\mathcal{F}_i(m)$ , the encoded map still corresponds to the observation of a single solution to a PDE,  $\mathcal{F}(\mathbf{w}^i; m) = Bu_i$ , albeit this time  $u_i$  solves the PDE  $A(m)u_i = f(\mathbf{w}^i)$ , i.e., with an encoded right-hand side.

### 3. Bayesian formulation of the inverse problem with encoded sources

This section contains a brief presentation of the Bayesian formulation of inverse problems with infinite-dimensional inversion parameters; for details we refer the reader to [19, 20] for theory and to [21] for the numerical approximation. In the Bayesian framework, the unknown parameter function  $m$  is modeled as a random field. Starting from a prior distribution law for  $m$ , we use observation data to obtain an improved description of the law of  $m$ . This updated distribution law of  $m$  is called the posterior measure. The prior measure, which we denote by  $\mu_0$ , can be understood as a probabilistic model for our prior beliefs about the parameter field  $m$ . The posterior measure, which we denote by  $\mu_{\text{post}}$ , is the distribution law of  $m$ , conditioned on observation data. A key ingredient of a Bayesian inverse problem is the data likelihood,  $\pi_{\text{like}}(\mathbf{d}|m)$ , which describes the conditional distribution of the data given the parameter field  $m$ ; this is where the parameter-to-observable map enters the Bayesian inverse problem.

Let  $\mathcal{D} \subset \mathbb{R}^d$  be a bounded domain with piecewise smooth boundary and  $(\Omega, \Sigma, \mathbb{P})$  a probability space. We consider an inference parameter  $m = m(x, \omega)$ , with  $(x, \omega) \in \mathcal{D} \times \Omega$ , such that for any  $\omega \in \Omega$ ,  $m(\cdot, \omega) \in \mathcal{V}$  where, as before,  $\mathcal{V}$  is an infinite-dimensional Hilbert space. Considering the law of  $m$  as a probability measure on  $(\mathcal{V}, \mathfrak{B}(\mathcal{V}))$ , the infinite-dimensional Bayes' theorem relates the Radon-Nikodym derivative of  $\mu_{\text{post}}$  with respect to  $\mu_0$  with the data likelihood  $\pi_{\text{like}}(\mathbf{d}|m)$ :

$$\frac{d\mu_{\text{post}}}{d\mu_0} \propto \pi_{\text{like}}(\mathbf{d}|m). \quad (4)$$

In the present work, we rely on Gaussian priors; i.e.,  $\mu_0 = \mathcal{N}(m_0, \mathcal{C}_0)$  is a Gaussian measure on  $\mathcal{V}$ . In that case, we require  $\mathcal{C}_0$  to be symmetric, positive and trace-class [19]. A common choice for  $\mathcal{C}_0$  (in two and three space dimensions) is the squared inverse of a Laplacian-like operator  $\mathcal{K}$ , i.e.,  $\mathcal{C}_0 = \mathcal{K}^{-2}$ . We also assume that the noise in the data is additive, and independent and identically distributed (over the different experiments); the distribution of each noise vector is normal with mean zero and covariance matrix  $\mathbf{\Gamma}_{\text{noise}}$ . That is,  $\mathbf{d}_i|m \sim \mathcal{N}(\mathcal{F}_i(m), \mathbf{\Gamma}_{\text{noise}})$ , for any  $i \in \{1, \dots, N_s\}$ . Consequently, each encoded observation  $\mathbf{d}(\mathbf{w}^i)$  will be normally distributed with mean zero and covariance matrix  $\mathbf{\Gamma}_{\text{noise},i} := (\sum_{j=1}^{N_s} (w_j^i)^2) \mathbf{\Gamma}_{\text{noise}}$ , i.e.,  $\mathbf{d}(\mathbf{w}^i)|m \sim \mathcal{N}(\mathcal{F}(\mathbf{w}^i; m), \mathbf{\Gamma}_{\text{noise},i})$ , for  $i \in \{1, \dots, N_w\}$ . Therefore, the likelihood function has the form

$$\pi_{\text{like}}(\mathbf{d}(\mathbf{w})|m) \propto \exp \left( -\frac{1}{2N_w} \sum_{i=1}^{N_w} \|\mathcal{F}(\mathbf{w}^i; m) - \mathbf{d}(\mathbf{w}^i)\|_{\mathbf{\Gamma}_{\text{noise},i}^{-1}}^2 \right).$$

### 3.1. MAP point

In finite dimensions, the MAP point is the parameter  $m_{\text{MAP}}$  that maximizes the posterior probability density function. Although this definition does not extend directly to the infinite-dimensional case, a MAP point can still be defined as a minimizer of a regularized data-misfit cost functional over an appropriate Hilbert subspace of the parameter space [19]. Let us define the Cameron-Martin space  $\mathcal{E} = \text{Im}(\mathcal{C}_0^{1/2})$ , endowed with the inner-product

$$\langle x, y \rangle_{\mathcal{E}} := \langle \mathcal{C}_0^{-1/2} x, \mathcal{C}_0^{-1/2} y \rangle = \langle \mathcal{K}x, \mathcal{K}y \rangle, \quad \forall x, y \in \mathcal{E}. \quad (5)$$

Then the MAP point is defined as

$$m_{\text{MAP}}(\mathbf{w}) = \arg \min_{m \in \mathcal{E}} \{ \mathcal{J}(\mathbf{w}; m) \}, \quad (6)$$

where, for the inverse problems considered in the present work, the functional  $\mathcal{J}(\mathbf{w}; \cdot) : \mathcal{E} \rightarrow \mathbb{R}$  is defined as

$$\mathcal{J}(\mathbf{w}; m) := \frac{1}{2N_w} \sum_{i=1}^{N_w} \|\mathcal{F}(\mathbf{w}^i; m) - \mathbf{d}(\mathbf{w}^i)\|_{\mathbf{\Gamma}_{\text{noise},i}^{-1}}^2 + \frac{1}{2} \|m - m_0\|_{\mathcal{E}}^2. \quad (7)$$

Here, the function  $m_0 \in \mathcal{E}$  is the mean of the prior measure.

### 3.2. Approximation to the posterior covariance

In general, there are no closed-form expressions for moments of the posterior measure. Thus, one usually relies on sampling-based methods to explore the posterior. For inverse problems governed by PDEs and problems with high-dimensional parameters (as, for instance, arising upon discretization of an infinite-dimensional parameter field), sampling of the posterior can quickly become infeasible since every evaluation of the likelihood requires a PDE solve. We thus rely on approximations of the posterior, namely Gaussian approximations about the MAP estimate. After finding the MAP point, we consider two commonly used approximations of the posterior measure by a Gaussian measure  $\mathcal{N}(m_{\text{MAP}}, \mathcal{C}_{\text{post}})$ , as discussed next [21, 22].

*Gauss-Newton approximation* Assuming the parameter-to-observable map  $\mathcal{F}(\mathbf{w}^i; \cdot)$  is Fréchet differentiable at the MAP point, one strategy to approximate the posterior is to linearize around the MAP point, i.e.,

$$\mathcal{F}(\mathbf{w}^i; m) \approx \mathcal{F}(\mathbf{w}^i; m_{\text{MAP}}) + \mathbf{J}_{\mathbf{w}^i}(m - m_{\text{MAP}}),$$

with  $\mathbf{J}_{\mathbf{w}^i} : \mathcal{V} \rightarrow \mathbb{R}$  the Fréchet derivative of the parameter-to-observable map  $\mathcal{F}(\mathbf{w}^i; \cdot)$  evaluated at the MAP point (6). Calling  $(\mathbf{J}_{\mathbf{w}^i})^*$  the adjoint of  $\mathbf{J}_{\mathbf{w}^i}$ , the covariance operator of the resulting Gaussian approximation of the posterior is given by

$$\mathcal{C}_{\text{post}}^{\text{GN}} = \left( \frac{1}{N_w} \sum_{i=1}^{N_w} (\mathbf{J}_{\mathbf{w}^i})^* \mathbf{\Gamma}_{\text{noise},i}^{-1} \mathbf{J}_{\mathbf{w}^i} + \mathcal{C}_0^{-1} \right)^{-1}. \quad (8)$$

Note that the operator that appears inside the brackets in (8) is the so called Gauss-Newton Hessian of the functional (7) evaluated at the MAP point,

$$\mathcal{H}_{\text{GN}}(m_{\text{MAP}}) := \frac{1}{N_w} \sum_{i=1}^{N_w} (\mathbf{J}_{\mathbf{w}^i})^* \mathbf{\Gamma}_{\text{noise},i}^{-1} \mathbf{J}_{\mathbf{w}^i} + \mathcal{C}_0^{-1}.$$

*Laplace approximation* Assuming  $\mathcal{J}(\mathbf{w}; \cdot)$ , in (7), is at least twice Fréchet differentiable at the MAP point, a second approach called Laplace approximation consists of using the second derivative of  $\mathcal{J}(\mathbf{w}; \cdot)$ , i.e., the Hessian, at the MAP point as an approximation to the posterior covariance

$$\mathcal{C}_{\text{post}}^{\text{L}} = (\mathcal{J}''(\mathbf{w}; m_{\text{MAP}}))^{-1} = \mathcal{H}^{-1}(m_{\text{MAP}}), \quad (9)$$

where the derivative in  $\mathcal{J}''$  is taken in terms of the parameter field  $m$ . Note that the Laplace approximation can be related, in finite dimensions, to a quadratic local approximation of  $\mathcal{J}(\mathbf{w}; \cdot)$  around the MAP point.

#### 4. A-optimal approach to source encoding

Combining the results from section 3 with elements from optimal experimental design, we propose a rigorous method to compute A-optimal encoding weights. In the Bayesian framework, the posterior covariance quantifies the uncertainty in the reconstruction. Since the posterior covariance depends on the weights (see section 4.1), we can select the weights that lead to a reconstruction with the least uncertainty. In the field of optimal experimental design, there are various design criteria that measure the statistical quality of the reconstructed parameter field [23]. In the present work, we rely on the A-optimal design criterion [23, 24], which aims to minimize the trace of the posterior covariance, or equivalently, to minimize the average posterior variance. That is, we compute the weights with the smallest trace of the posterior covariance  $\Phi(\mathbf{w}) = \text{tr}(\mathcal{C}_{\text{post}})$ , with  $\mathcal{C}_{\text{post}}$  given by  $\mathcal{C}_{\text{post}}^{\text{GN}}$  (8) or  $\mathcal{C}_{\text{post}}^{\text{L}}$  (9).

An alternate view of the A-optimal design criterion is that of minimizing the expected Bayes risk of the MAP estimator, which coincides with the trace of the posterior covariance for a linear inverse problem [9, 11, 25]. This interpretation of the A-optimal criterion can be stated as the average mean squared error between the MAP estimator (i.e., the parameter reconstruction) and the true parameter (e.g., see [9]). While this interpretation of A-optimality is restricted to linear inverse



problems, it provides another motivation for our choice of the design criterion. In our numerical results, we explore this relation between minimizing the trace of the posterior covariance and the mean squared distance between the MAP point and the true parameter and observe that minimizing the trace of the posterior covariance correlates with smaller errors for the parameter reconstruction.

#### 4.1. Dependence of the operators $\mathcal{C}_{\text{post}}^{\text{GN}}$ and $\mathcal{C}_{\text{post}}^{\text{L}}$ on $\mathbf{w}$

The dependence of the operators  $\mathcal{C}_{\text{post}}^{\text{GN}}$  (8) and  $\mathcal{C}_{\text{post}}^{\text{L}}$  (9) on the weights is twofold. First these operators depend on the encoded parameter-to-observable maps that depend explicitly on the weights,  $\mathcal{F}(\mathbf{w}^i; m) = \sum_{j=1}^{N_s} w_j^i \mathcal{F}_j(m)$ . Moreover, the posterior covariance operators also depend on the weights through the MAP point (6), which depends on the weights as illustrated by (6) and (7).

The dependence of the covariance operator  $\mathcal{C}_{\text{post}}^{\text{GN}}$  on  $\mathbf{w}$  is straightforward to see. In particular, using the chain-rule on the forward problem  $\mathcal{A}(m)u_i = f(\mathbf{w}^i)$ , the Fréchet derivative of the parameter-to-observable at the MAP point is given by

$$\mathbf{J}_{\mathbf{w}^i} = -B\mathcal{A}(m_{\text{MAP}}(\mathbf{w}))^{-1} \frac{\partial \mathcal{A}(m)u_i}{\partial m} \bigg|_{m=m_{\text{MAP}}(\mathbf{w})}. \quad (10)$$

Given  $N_w$  encoding weights  $\mathbf{w} = (\mathbf{w}^1, \dots, \mathbf{w}^{N_w})$  where  $\mathbf{w}^i \in \mathbb{R}^{N_s}$ , we emphasize the dependence of the posterior covariance on the weights by writing  $\mathcal{C}_{\text{post}}^{\text{GN}} = \mathcal{C}_{\text{post}}^{\text{GN}}(\mathbf{w})$ . The structure of the covariance operator  $\mathcal{C}_{\text{post}}^{\text{L}}$  is more complicated. We detail the dependence of  $\mathcal{C}_{\text{post}}^{\text{L}}$  on  $\mathbf{w}$  for the application problem considered in the present paper in section 5. Note that in the case of a linear parameter-to-observable map, both posterior covariances (8) and (9) are equal.

In the present formulation,  $\text{tr}(\mathcal{C}_{\text{post}}(\mathbf{w}))$  scales with the weights. For instance, applying a constant multiplicative factor  $\lambda > 1$  to all weights would reduce the influence of the prior in the computation of the MAP point (6) for once. It would also inflate the norm of the state variable  $u_i$  by that factor  $\lambda$ , which would then increase the size of the derivative (10). This would in turn artificially reduce the trace of the posterior covariance (8). A solution is to restrict the codomain of each encoding weight to a sphere of radius  $r$  in  $\mathbb{R}^{N_s}$ . We denote the corresponding space, for the weights  $\mathbf{w}$ , by  $\mathcal{S}_r$ , i.e.,  $\mathcal{S}_r := \{\mathbf{w} = (\mathbf{w}^1, \dots, \mathbf{w}^{N_w}) \in \mathbb{R}^{N_w N_s}; |\mathbf{w}^i| = r, \forall i\}$ . As discussed in section 2, the theory of randomized trace estimation dictates the use of  $r = \sqrt{N_s}$ . However this value is arbitrary and can be compensated by an equivalent re-scaling of the regularization parameter. Therefore for simplicity we use the value  $r = 1$  along with the notation  $\mathcal{S} := \mathcal{S}_1$ . Another implication of that choice,  $|\mathbf{w}^i| = 1$ , is that the covariance matrices for the encoded noise vectors, introduced in section 3, simplify to  $\mathbf{\Gamma}_{\text{noise},i} = \mathbf{\Gamma}_{\text{noise}}$ , for  $i \in \{1, \dots, N_w\}$ .

#### 4.2. A-optimal encoding weights

We propose to compute the A-optimal encoding weights as the solution to the constrained minimization problem

$$\min_{\mathbf{w} \in \mathcal{S}} \Phi(\mathbf{w}) := \text{tr}(\mathcal{C}_{\text{post}}(\mathbf{w})). \quad (11)$$

Since there are no closed-form expressions for moments of the posterior measure, we replace the exact posterior covariance in (11) with one of the two approximations



introduced in section 3.2. The Gauss–Newton formulation of the A-optimal encoding weights,

$$\Phi_{\text{GN}}(\mathbf{w}) = \text{tr}(\mathcal{H}_{\text{GN}}^{-1}(\mathbf{w}; m_{\text{MAP}}(\mathbf{w}))), \quad (12)$$

is based on the posterior covariance approximation (8), and the Laplace formulation,

$$\Phi_{\text{L}}(\mathbf{w}) = \text{tr}(\mathcal{H}^{-1}(\mathbf{w}; m_{\text{MAP}}(\mathbf{w}))), \quad (13)$$

is based on the posterior covariance (9). Note that both formulations (12) and (13) require the computation of the MAP point which is computationally expensive for large-scale problems. To avoid the cost associated with the computation of the MAP point, an additional simplification of (12) can be achieved by evaluating the posterior covariance (8) at a reference parameter field  $m_0$ , which leads to the following (simplified) objective function,

$$\Phi_0(\mathbf{w}) = \text{tr}(\mathcal{H}_{\text{GN}}^{-1}(\mathbf{w}; m_0)). \quad (14)$$

*A-optimal encoding weights formulation for large-scale applications* Formulation (11) is a nonlinear optimization problem that requires the use of iterative methods. These methods involve repeated evaluations of the trace of the posterior covariance. Following discretization, the posterior covariance is a high-dimensional operator that is defined implicitly, i.e., through its applications to vectors. The exact computation of the trace of such operators, and their derivatives with respect to encoding weights, is computationally intractable. For this reason, we propose an approximate formulation using a randomized trace estimator (see [15, 16] for the theory, and [8, 9] for examples of applications). Following the formulation in [10], we introduce the Gaussian measure  $\mu_\delta = \mathcal{N}(0, \mathcal{C}_\delta)$  where  $\mathcal{C}_\delta := (I - \delta\Delta)^{-2}$ . Here  $\Delta$  denotes the Laplacian operator with homogeneous Neumann boundary conditions and  $\delta > 0$  a sufficiently small real number. Then for any positive, self-adjoint and trace-class operator  $\mathcal{T}$ , we may use an estimator of form,

$$\text{tr}(\mathcal{T}) \approx \frac{1}{n_{tr}} \sum_{i=1}^{n_{tr}} \langle \mathcal{T} z_i, z_i \rangle_{\mathcal{H}},$$

where the  $z_i$  are drawn from  $\mu_\delta$ . In practice, reasonable approximations of the trace can be obtained with a relatively small  $n_{tr}$ .

The optimization problem for finding A-optimal encoding weights is formulated as follows

$$\min_{\mathbf{w} \in \mathcal{S}} \frac{1}{n_{tr}} \sum_{i=1}^{n_{tr}} \langle \mathcal{C}_{\text{post}}(\mathbf{w}) z_i, z_i \rangle.$$

Specializing to the cases of  $\Phi_{\text{GN}}(\mathbf{w})$  (12) and  $\Phi_{\text{L}}(\mathbf{w})$  (13) results in the following formulations,

$$\min_{\mathbf{w} \in \mathcal{S}} \left\{ \frac{1}{n_{tr}} \sum_{i=1}^{n_{tr}} \langle \mathcal{H}_{\text{GN}}^{-1}(\mathbf{w}; m_{\text{MAP}}(\mathbf{w})) z_i, z_i \rangle \right\}, \quad (15)$$

$$\min_{\mathbf{w} \in \mathcal{S}} \left\{ \frac{1}{n_{tr}} \sum_{i=1}^{n_{tr}} \langle \mathcal{H}^{-1}(\mathbf{w}; m_{\text{MAP}}(\mathbf{w})) z_i, z_i \rangle \right\}. \quad (16)$$

Again to avoid the cost associated with the computation of the MAP point, one can evaluate the Gauss–Newton Hessian in (15) at a fixed reference parameter field  $m_0$ ; this leads to the following (simplified) optimization problem,

$$\min_{\mathbf{w} \in \mathcal{S}} \left\{ \frac{1}{n_{tr}} \sum_{i=1}^{n_{tr}} \langle \mathcal{H}_{\text{GN}}^{-1}(\mathbf{w}; m_0) z_i, z_i \rangle \right\}. \quad (17)$$

The formulation (17) can be seen as an extension of the formulation proposed in [8] to a fully nonlinear inverse problem formulated at the infinite-dimensional level.

## 5. Application to the Helmholtz inverse problem

In this section, we elaborate the A-optimal encoding weights formulation introduced in section 4 for the Helmholtz inverse problem. Recall that high resolution reconstructions in this application require a large number of experiments and that the computational cost of the inversion scales linearly with the number of experiments (see section 1). Source encoding can provide a trade-off between high-quality reconstruction and computational cost.

We begin by describing the inverse problem used in our study (section 5.1). Then the optimization problem to compute the A-optimal encoding weights, including the adjoint-based expressions for the gradient of this objective function, is detailed in section 5.2.

### 5.1. The inverse problem: medium parameter reconstruction

For simplicity of the presentation, we derive the formulation using a single frequency but extensions to the case of multiple frequencies are straightforward. We use homogeneous Neumann boundary conditions. The frequency-domain Helmholtz equation is given, for  $i = 1, \dots, N_w$ , by

$$\begin{aligned} -\Delta u_i - \kappa^2 m u_i &= f(\mathbf{w}^i) && \text{in } \mathcal{D}, \\ \nabla u_i \cdot \mathbf{n} &= 0 && \text{on } \partial\mathcal{D}. \end{aligned} \quad (18)$$

Solutions  $u_i$  (18) are considered in  $H^1(\mathcal{D})$ , i.e., the Sobolev space of functions in  $L^2(\mathcal{D})$  with square integrable weak derivatives. The original source terms are in the dual space of  $H_0^1(\mathcal{D})$ , i.e.,  $f_j \in H^{-1}(\mathcal{D})$ . The (medium) parameter field  $m \in L^\infty(\mathcal{D})$  corresponds to the square of the slowness (or the squared inverse local wave speed) and the constant  $\kappa$  is the frequency of the wave (in rad/s).

**5.1.1. MAP point** The MAP point is the solution to a deterministic inverse problem (see section 3.1) with the norms in the data-misfit and regularization terms weighted by the noise and prior covariance operators respectively. In particular, with a Gaussian prior  $\mu_0 = \mathcal{N}(m_0, \mathcal{C}_0)$  and the norm corresponding to the inner product (5), we have

$$m_{\text{MAP}}(\mathbf{w}) = \arg \min_{m \in \mathcal{E}} \left\{ \frac{1}{2N_w} \sum_{i=1}^{N_w} \|B u_i - \mathbf{d}(\mathbf{w}^i)\|_{\Gamma_{\text{noise}}^{-1}}^2 + \frac{1}{2} \|m - m_0\|_{\mathcal{E}}^2 \right\}, \quad (19)$$

where  $u_i$  solves (18).

To properly define the source terms  $f_i$ , appearing in the right hand-side of the forward problem, and the observation operator  $B$ , we define the mollifier  $\varphi_\varepsilon(x; y)$  as follows:

$$\varphi_\varepsilon(x; y) = \frac{1}{\alpha_\varepsilon} e^{-\frac{1}{\varepsilon^2 - |x-y|^2}} \mathbf{1}_{\mathcal{B}(y, \varepsilon)}(x), \quad (20)$$

where  $\alpha_\varepsilon = 2\pi K \varepsilon^2 e^{-1/\varepsilon^2}$ ,  $K = \int_0^1 r e^{-1/(1-r^2)} dr$ ,  $\mathbf{1}_{\mathcal{B}(y, \varepsilon)}$  is the indicator function for the ball of radius  $\varepsilon$  centered at  $y$ , and  $0 < \varepsilon \ll 1$ . This function is smooth and integrates to one. We choose each source terms  $f_i$  to be a mollifier centered at one of the  $N_s$  source locations that we denote  $x_i^s$  for  $i = 1, \dots, N_s$ , i.e.,  $f_i(x) = \varphi_\varepsilon(x; x_i^s)$ . The observation operator  $B : H^1(\mathcal{D}) \rightarrow \mathbb{R}^q$  is the evaluation, at each of the receiver locations which we denote  $x_j^r$  for  $j = 1, \dots, q$ , of a convolution between the solution to the forward problem  $u_i$  and a mollifier  $\varphi_{\varepsilon'}(x; 0)$ , i.e.,  $(Bu_i)_j = (u_i * \varphi_{\varepsilon'}(\cdot; 0))(x_j^r)$ . These choices of the source terms and observation operator guarantee that the forward, adjoint, incremental forward and incremental adjoint solutions belong to  $H^1(\mathcal{D})$ .

*5.1.2. Gradient and Hessian of the inverse problem* Availability of derivatives of the function in brackets on the right hand side of (19) is required for the computation of  $m_{\text{MAP}}$ . The second derivative, i.e., the Hessian operator, also enters the A-optimal formulation laid down in section 4. We derive both gradient and Hessian following the formal Lagrangian approach [1, 26]. The first-order necessary optimality condition for the MAP point is a coupled system of PDEs: Find  $(m_{\text{MAP}}, \{u_i\}_i, \{p_i\}_i) \in \mathcal{E} \times H^1(\mathcal{D})^{N_w} \times H^1(\mathcal{D})^{N_w}$  such that for all variations  $(\tilde{m}, \{\tilde{u}_i\}_i, \{\tilde{p}_i\}_i) \in \mathcal{E} \times H^1(\mathcal{D})^{N_w} \times H^1(\mathcal{D})^{N_w}$

$$\begin{aligned} \langle \nabla u_i, \nabla \tilde{p}_i \rangle - \kappa^2 \langle m_{\text{MAP}}(\mathbf{w}) u_i, \tilde{p}_i \rangle - \langle f(\mathbf{w}^i), \tilde{p}_i \rangle &= 0, \forall i \\ \langle \nabla \tilde{u}_i, \nabla p_i \rangle - \kappa^2 \langle \tilde{u}_i, m_{\text{MAP}}(\mathbf{w}) p_i \rangle + \langle B \tilde{u}_i, Bu_i - \mathbf{d}(\mathbf{w}^i) \rangle_{\mathbf{r}_{\text{noise}}^{-1}} &= 0, \forall i \\ \langle m_{\text{MAP}}(\mathbf{w}) - m_0, \tilde{m} \rangle_{\mathcal{E}} - \frac{1}{N_w} \sum_{i=1}^{N_w} \kappa^2 \langle u_i p_i, \tilde{m} \rangle &= 0. \end{aligned} \quad (21)$$

For the Hessian, we describe the solution to the equation  $y = \mathcal{H}^{-1}(m_{\text{MAP}})z$ . This leads to the coupled system of PDEs: Find  $(y, \{v_i\}_i, \{q_i\}_i) \in \mathcal{E} \times H^1(\mathcal{D})^{N_w} \times H^1(\mathcal{D})^{N_w}$  such that for all  $(\tilde{m}, \{\tilde{u}_i\}_i, \{\tilde{p}_i\}_i) \in \mathcal{E} \times H^1(\mathcal{D})^{N_w} \times H^1(\mathcal{D})^{N_w}$  the following equations are satisfied:

$$\begin{aligned} \langle \nabla v_i, \nabla \tilde{p}_i \rangle - \kappa^2 \langle m_{\text{MAP}}(\mathbf{w}) v_i, \tilde{p}_i \rangle - \kappa^2 \langle u_i y, \tilde{p}_i \rangle &= 0, \forall i \\ \langle \nabla \tilde{u}_i, \nabla q_i \rangle - \kappa^2 \langle \tilde{u}_i, m_{\text{MAP}}(\mathbf{w}) q_i \rangle - \kappa^2 \langle \tilde{u}_i, p_i y \rangle + \langle B \tilde{u}_i, B v_i \rangle_{\mathbf{r}_{\text{noise}}^{-1}} &= 0, \forall i \\ \langle y, \tilde{m} \rangle_{\mathcal{E}} - \frac{1}{N_w} \sum_{i=1}^{N_w} \kappa^2 \left[ \langle v_i p_i, \tilde{m} \rangle + \langle u_i q_i, \tilde{m} \rangle \right] &= \langle z, \tilde{m} \rangle. \end{aligned} \quad (22)$$

## 5.2. The optimization problem for A-optimal encoding weights

Here we formulate the optimization problem for computing A-optimal source encoding weights for the frequency-domain seismic inverse problem (18). We restrict ourselves to the case of the Laplace formulation (16) as the other two functionals, (15) and (17), can be treated as special cases of the Laplace formulation.

In its original format, the optimization problem for A-optimal encoding weights (16) is a bi-level optimization, as the MAP point is itself the solution to

a minimization problem (6). However this is not a practical formulation to compute derivatives. We therefore reformulate (16) as a PDE-constrained optimization problem in which the MAP point is defined as a solution of the first-order optimality condition (21). The other PDE constraint is the solution to the Hessian system (22) along the random directions of the trace estimator, i.e., we define the objective functional for the computation of the A-optimal encoding weights by

$$\frac{1}{n_{tr}} \sum_{k=1}^{n_{tr}} \langle y_k, z_k \rangle,$$

where  $z_k$  is a random direction for the trace estimator and  $y_k = \mathcal{H}^{-1}(m_{\text{MAP}})z_k$  according to (22). We can then enforce these PDE constraints with Lagrange multipliers and compute derivatives of the optimization problem (16) using the formal Lagrangian approach. We account for the constraint on the weights through a penalty term,

$$\frac{\lambda}{2N_w} \sum_{j=1}^{N_w} (\|\mathbf{w}^j\|^2 - 1)^2,$$

with  $\lambda \in \mathbb{R}$ . Although a penalty term is not the only option, we found this relaxation of the constraint to be efficient and easy to implement.

We now present the complete formulation for (16). The A-optimal encoding weights are solutions to the minimization problem

$$\min_{\mathbf{w}} \left\{ \frac{1}{n_{tr}} \sum_{k=1}^{n_{tr}} \langle y_k, z_k \rangle + \frac{\lambda}{2N_w} \sum_{j=1}^{N_w} (\|\mathbf{w}^j\|^2 - 1)^2 \right\}, \quad (23)$$

where for every  $k = 1, \dots, n_{tr}$ ,  $(y_k, \{v_{i,k}\}_i, \{q_{i,k}\}_i) \in \mathcal{E} \times H^1(\mathcal{D})^{N_w} \times H^1(\mathcal{D})^{N_w}$  solves the system

$$\begin{aligned} \langle \nabla v_{i,k}, \nabla \tilde{p}_{i,k} \rangle - \kappa^2 \langle m_{\text{MAP}}(\mathbf{w}) v_{i,k}, \tilde{p}_{i,k} \rangle - \kappa^2 \langle u_i y_k, \tilde{p}_{i,k} \rangle &= 0, \forall i \\ \langle \nabla \tilde{u}_{i,k}, \nabla q_{i,k} \rangle - \kappa^2 \langle \tilde{u}_{i,k}, m_{\text{MAP}}(\mathbf{w}) q_{i,k} \rangle - \kappa^2 \langle \tilde{u}_{i,k}, p_i y_k \rangle \\ &\quad + \langle B \tilde{u}_{i,k}, B v_{i,k} \rangle_{\mathbf{\Gamma}_{\text{noise}}^{-1}} = 0, \forall i \\ \langle y_k, \tilde{m} \rangle_{\mathcal{E}} - \frac{1}{N_w} \sum_{i=1}^{N_w} \kappa^2 \left[ \langle v_{i,k} p_i, \tilde{m} \rangle + \langle u_i q_{i,k}, \tilde{m} \rangle \right] &= \langle z_k, \tilde{m} \rangle, \end{aligned} \quad (24)$$

for all  $(\tilde{m}, \{\tilde{u}_{i,k}\}_i, \{\tilde{p}_{i,k}\}_i) \in \mathcal{E} \times H^1(\mathcal{D})^{N_w} \times H^1(\mathcal{D})^{N_w}$  and where  $(m_{\text{MAP}}, \{u_i\}_i, \{p_i\}_i) \in \mathcal{E} \times H^1(\mathcal{D})^{N_w} \times H^1(\mathcal{D})^{N_w}$  solves the first-order optimality system for the Helmholtz inverse problem

$$\begin{aligned} \langle \nabla u_i, \nabla \tilde{p}_i \rangle - \kappa^2 \langle m_{\text{MAP}}(\mathbf{w}) u_i, \tilde{p}_i \rangle - \langle f(\mathbf{w}^i), \tilde{p}_i \rangle &= 0, \forall i \\ \langle \nabla \tilde{u}_i, \nabla p_i \rangle - \kappa^2 \langle \tilde{u}_i, m_{\text{MAP}}(\mathbf{w}) p_i \rangle + \langle B \tilde{u}_i, B u_i - \mathbf{d}(\mathbf{w}^i) \rangle_{\mathbf{\Gamma}_{\text{noise}}^{-1}} &= 0, \forall i \\ \langle m_{\text{MAP}}(\mathbf{w}) - m_0, \tilde{m} \rangle_{\mathcal{E}} - \frac{1}{N_w} \sum_{i=1}^{N_w} \kappa^2 \langle u_i p_i, \tilde{m} \rangle &= 0, \end{aligned}$$

for all  $(\tilde{m}, \{\tilde{u}_i\}_i, \{\tilde{p}_i\}_i) \in \mathcal{E} \times H^1(\mathcal{D})^{N_w} \times H^1(\mathcal{D})^{N_w}$ .

**5.2.1. Gradient of the A-optimal weight problem** We derive the gradient of the objective function defined in (23), with respect to  $\mathbf{w}$ , using a formal Lagrangian approach. We refer the reader to [Appendix A](#) for this derivation. Since we enforce the PDE constraints weakly using Lagrange multipliers, we introduce adjoint variables that are indicated with a star superscript, e.g.,  $m^*$  is the adjoint variable for  $m$ . The gradient is given by  $[\delta_{\mathbf{w}^1} \Phi_L(\mathbf{w}), \delta_{\mathbf{w}^2} \Phi_L(\mathbf{w}), \dots, \delta_{\mathbf{w}^{N_w}} \Phi_L(\mathbf{w})]^T$ , where for any  $i = 1, \dots, N_w$ ,

$$\delta_{\mathbf{w}^i} \Phi_L(\mathbf{w}) = -\frac{1}{N_w} \begin{bmatrix} \langle f_1, u_i^* \rangle + \langle Bp_i^*, \mathbf{d}_1 \rangle_{\Gamma_{\text{noise}}^{-1}} \\ \langle f_2, u_i^* \rangle + \langle Bp_i^*, \mathbf{d}_2 \rangle_{\Gamma_{\text{noise}}^{-1}} \\ \vdots \\ \langle f_{N_s}, u_i^* \rangle + \langle Bp_i^*, \mathbf{d}_{N_s} \rangle_{\Gamma_{\text{noise}}^{-1}} \end{bmatrix}.$$

The variables  $u_i^*$  and  $p_i^*$  are computed by solving the following Hessian-like system (compare with (22)): Find  $(m^*, \{u_i^*\}_i, \{p_i^*\}_i) \in \mathcal{E} \times H^1(\mathcal{D})^{N_w} \times H^1(\mathcal{D})^{N_w}$  such that for all  $(\tilde{m}, \{\tilde{u}_i\}_i, \{\tilde{p}_i\}_i) \in \mathcal{E} \times H^1(\mathcal{D})^{N_w} \times H^1(\mathcal{D})^{N_w}$  the following equations are satisfied:

$$\begin{aligned} \langle \nabla p_i^*, \nabla \tilde{p}_i \rangle - \kappa^2 \langle mp_i^*, \tilde{p}_i \rangle - \kappa^2 \langle u_i m^*, \tilde{p}_i \rangle &= -\frac{2}{n_{tr}} \sum_{k=1}^{n_{tr}} \kappa^2 \langle v_{i,k} y_k, \tilde{p}_i \rangle, \\ \langle \nabla u_i^*, \nabla \tilde{u}_i \rangle - \kappa^2 \langle mu_i^*, \tilde{u}_i \rangle - \kappa^2 \langle p_i m^*, \tilde{u}_i \rangle \\ &\quad + \langle Bp_i^*, B\tilde{u}_i \rangle_{\Gamma_{\text{noise}}^{-1}} = -\frac{2}{n_{tr}} \sum_{k=1}^{n_{tr}} \kappa^2 \langle y_k q_{i,k}, \tilde{u}_i \rangle, \\ \langle m^*, \tilde{m} \rangle_{\mathcal{E}} - \frac{1}{N_w} \sum_{i=1}^{N_w} \kappa^2 [\langle u_i u_i^*, \tilde{m} \rangle + \langle p_i^* p_i, \tilde{m} \rangle] &= -\frac{2}{n_{tr} N_w} \sum_{k=1}^{n_{tr}} \sum_{i=1}^{N_w} \kappa^2 \langle v_{i,k} q_{i,k}, \tilde{m} \rangle. \end{aligned} \tag{25}$$

The variables  $\{v_{i,k}\}$  (resp.  $\{q_{i,k}\}$ ) are the incremental state (resp. adjoint) variables which occur in the application of the inverse Hessian in the direction of the  $k$ -th trace estimator  $z_k$ .

**5.2.2. Discretization** The numerical solution of (23) is done via the Optimize-then-Discretize (OTD) approach, where the discretization is based on continuous Galerkin finite element with Lagrange nodal basis functions. Extra care is needed for the discretization of the covariance operator to ensure that its discrete representation faithfully represents the properties of the target infinite-dimensional object. We do not provide full details of the discretization and refer the reader to [9, 21]. However, we show how to select the discrete random directions  $z_k$  in the trace estimator. Let us call  $V_h$  the finite-dimensional approximation to the space  $H^1(\mathcal{D})$  used for the finite-element representations of all state, adjoint, corresponding incremental variables and their respective adjoints. And let  $V_h^m$  be the finite-dimensional space for the medium parameter  $m$ . Let us call  $\{\psi_i\}_{i=1}^l$  (resp.  $\{\phi_i\}_{i=1}^l$ ) a basis for  $V_h$  (resp.  $V_h^m$ ). Let us introduce the vector notations  $\mathbf{y}_k = (y_k^1, \dots, y_k^l)^T$  (resp.  $\mathbf{z}_k = (z_k^1, \dots, z_k^l)^T$ ) for the finite element representations of  $y_k$  (resp.  $z_k$ ) in  $V_h^m$ . The finite-dimensional approximation to the trace estimation is then

$$\frac{1}{n_{tr}} \sum_{k=1}^{n_{tr}} \langle y_k^h, z_k^h \rangle_{L^2} = \frac{1}{n_{tr}} \sum_{k=1}^{n_{tr}} \sum_{i,j=1}^l y_k^i z_k^j \langle \phi_i, \phi_j \rangle_{L^2} = \frac{1}{n_{tr}} \sum_{k=1}^{n_{tr}} \langle \mathbf{y}_k, \mathbf{z}_k \rangle_{\mathbf{M}},$$

with  $\mathbf{M}_{ij} = \langle \phi_i, \phi_j \rangle_{L^2}$  the mass matrix in  $V_h^m$ . From the definition of  $y_k$ , we see that each  $y_k^h$  solves the system  $\langle \mathcal{H}y_k^h, \phi_i \rangle_{L^2} = \langle z_k^h, \phi_i \rangle_{L^2}$ , for  $i = 1, \dots, l$ . Substituting the representation of  $y_k^h$  and  $z_k^h$  in the basis of  $V_h^m$ , we obtain the matrix system  $\mathbf{H}\mathbf{y}_k = \mathbf{M}\mathbf{z}_k$ , where  $\mathbf{H}$  is the standard Hessian matrix obtained from finite-element discretization of system (22), i.e.,  $\mathbf{H}_{ij} = \langle \mathcal{H}\phi_j, \phi_i \rangle_{L^2}$ . The finite-dimensional approximation to the trace estimation becomes

$$\frac{1}{n_{tr}} \sum_{k=1}^{n_{tr}} \langle y_k^h, z_k^h \rangle_{L^2} = \frac{1}{n_{tr}} \sum_{k=1}^{n_{tr}} \langle \mathbf{H}^{-1} \mathbf{M} \mathbf{z}_k, \mathbf{z}_k \rangle_{\mathbf{M}} = \frac{1}{n_{tr}} \sum_{k=1}^{n_{tr}} \langle \mathbf{H}_{\mathbf{M}}^{-1} \mathbf{z}_k, \mathbf{z}_k \rangle_{\mathbf{M}},$$

where we defined  $\mathbf{H}_{\mathbf{M}}^{-1} := \mathbf{H}^{-1} \mathbf{M}$ . The matrix  $\mathbf{H}_{\mathbf{M}}^{-1}$  is  $\mathbf{M}$ -symmetric [21], i.e., self-adjoint with respect to the  $\mathbf{M}$  inner-product. Then it was proved in [9] that  $\frac{1}{n_{tr}} \sum_{k=1}^{n_{tr}} \langle \mathbf{H}_{\mathbf{M}}^{-1} \mathbf{z}_k, \mathbf{z}_k \rangle_{\mathbf{M}}$  is indeed a trace estimator provided  $\mathbf{z}_k \sim \mathcal{N}(0, \mathbf{M}^{-1})$ . In practice, vectors  $\mathbf{z}_k$  are sampled by taking draws  $\mathbf{x}_k$  from multivariate standard normal distribution,  $\mathbf{x}_k \sim \mathcal{N}(0, \mathbf{I})$ , and using  $\mathbf{z}_k = \mathbf{M}^{-1/2} \mathbf{x}_k$ .

**5.2.3. Computational cost** Problem (23) is highly nonlinear and requires iterative methods to be solved. The gradient, derived in section 5.2.1, allows us to use quasi-Newton methods [27]. In table 1, we report the dominating terms of the computational cost of evaluating the objective function and its gradient in all three cases (15)–(17). Additionally, it is possible to reduce the cost of formulation (15) by computing a low-rank approximation of the Hessian operator [28]. One must keep in mind, however, that the incremental state variables  $\{v_{i,k}\}$  and incremental adjoint variables  $\{q_{i,k}\}$  corresponding to each random directions  $\{z_k\}$  are required to compute the gradient. For this reason, a low-rank approximation of the Hessian will only lower the computational cost when  $n_{tr} > n_{cg} n_{newt}$ .

Table 1: Computational cost for objective function and gradient evaluation of the optimization problem for finding A-optimal encoding weights. We report the computational cost, in terms of the number of forward PDE solves, for  $\Phi_{\text{GN}}(\mathbf{w})$ ,  $\Phi_{\text{L}}(\mathbf{w})$ , and  $\Phi_0(\mathbf{w})$  defined in (15)–(17) respectively. Notations:  $n_{\text{cg}}$  = number of Conjugate-Gradient iterations to compute the search direction in Newton’s method;  $n_{\text{newt}}$  = number of Newton steps to compute the MAP point.

	$\Phi_0(\mathbf{w})$	$\Phi_{\text{GN}}(\mathbf{w})$ and $\Phi_{\text{L}}(\mathbf{w})$ (no low-rank)	$\Phi_{\text{GN}}(\mathbf{w})$ (with low-rank)
objective evaluation			
MAP point	$2N_w$	$2N_w n_{\text{cg}} n_{\text{newt}}$	$2N_w n_{\text{cg}} n_{\text{newt}}$
$\text{tr}(\mathcal{H}^{-1})$	$2N_w n_{\text{cg}} n_{tr}$	$2N_w n_{\text{cg}} n_{tr}$	$2N_w n_{\text{cg}}$
gradient evaluation			
$v_{ik}, q_{ik}$	–	–	$2N_w n_{tr}$
$m^*$	–	$2N_w n_{\text{cg}}$	–
$u_i^*, p_i^*$	$N_w$	–	$2N_w$
total	$2N_w n_{\text{cg}} n_{tr}$	$2N_w n_{\text{cg}} (n_{\text{newt}} + n_{tr})$	$2N_w (n_{\text{cg}} n_{\text{newt}} + n_{tr})$

Following the OTD approach, the optimization problem (23) is formulated in function space, before being solved with algorithms that are discretization-independent. This results in the overall computational cost being independent of the discretization of the parameter space, or in other words, each of the quantities

$n_{\text{newt}}$ ,  $n_{\text{cg}}$  and  $n_{\text{tr}}$  in table 1 remain constant when the mesh gets refined. We spend the rest of this section discussing the choice of such discretization-invariant algorithms. First, we use Newton’s method, with Armijo line search, to compute the MAP point; the number of Newton steps needed to converge,  $n_{\text{newt}}$ , is typically independent of the size of the parameter space [29]. Moreover, the Hessian system (22) needed to compute the MAP point, to evaluate the objective function (23), and to compute the adjoint variable  $m^*$  (25), is solved using the preconditioned Conjugate Gradient method [27]. The Conjugate Gradient solver is preconditioned by the prior covariance operator; the number of iterations  $n_{\text{cg}}$  needed to solve the Hessian system then depends on the spectral properties of the prior-preconditioned data-misfit part of the Hessian operator (i.e., the Hessian in function space) and is therefore independent of the discretization. The trace estimator displays a similar type of behaviour. The number of trace estimator vectors  $n_{\text{tr}}$  one should use depends on the spectral properties of the underlying infinite-dimensional operator. The choice of a discrete inner-product weighted by the mass matrix (see section 5.2.2) guarantees that our discrete operator will be a valid approximation of the infinite-dimensional operator and will conserve its spectral properties. The actual evaluation of the trace is performed through the repeated solution of the Hessian system (24), which was shown above to be discretization-independent.

## 6. Numerical results

In this section, we present numerical results for the Helmholtz inverse problem in two (spatial) dimensions. We start with a low-dimensional example ( $N_w = 1$  for  $N_s = 2$ ), which allows us to visualize the objective functions defined in section 4.2 over the entire weight space. This facilitates a qualitative comparison of the different approximations introduced, the Gauss–Newton (12) and Laplace objective functions (13), along with the linearized formulation (14). We then present an example with a higher-dimensional weight space ( $N_s = 10$ ) in which we study the distribution of the A-optimal encoding weights and random weights sampled from the uniform spherical distribution and how the number of encoded weight vectors influence these results.

The setting for this section is a square domain with 20 receivers located at the top of the domain, and sources positioned on the bottom and left edges of the domain. The source term is a mollifier (20) with  $\varepsilon = 10^{-6}$ . This choice of source terms was numerically found to be reasonably well approximated, at the discrete level, by a point source; we utilize that approximation in this section. We use a wave frequency of  $\kappa = 2\pi$  in equation (18). All partial differential equations are discretized by continuous Galerkin finite elements (linear elements for the parameters and quadratic elements for the state and adjoint variables). This results in a (medium) parameter space of 182 degrees of freedom. We work with synthetic data that are polluted by a 2% additive Gaussian noise.

### 6.1. One-dimensional weight space

In this section, we study a one-dimensional source encoding problem corresponding to a single linear combination of two sources ( $N_s = 2$  and  $N_w = 1$ ). Although this setting represents an unrealistic situation (low number of sources, and high ratio of number of encoded sources over total number of sources), it is informative for the following reasons: (1) It provides numerical evidence of the strong and highly nonlinear



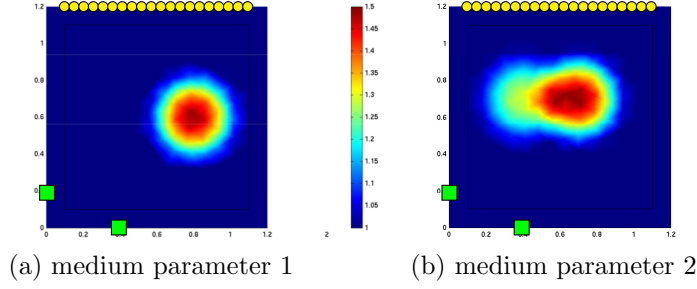


Figure 1: Target medium parameters, along with the locations of the sources (green squares) and receivers (yellow circles).

dependence of the objective functions (12)–(14) on the encoding weights. (2) It demonstrates the presence of multiple local minima in the minimization problem (11). (3) It highlights the difference between the Gauss–Newton and Laplace formulations. The sources are located on the bottom and left edges of the domain, and we study two different medium parameters, each made of a constant background and a smooth compactly supported perturbation (see figure 1).

We next define the noise covariance and the prior covariance operators used in these numerical applications. Let us introduce the non-singular, positive definite, elliptic operator  $\mathcal{Y} = -\gamma\Delta + \beta I$ , with  $\gamma, \beta$  positive constants,  $I$  the identity operator and  $\Delta$  the Laplacian operator with homogeneous Neumann boundary conditions. Then we define the prior covariance operator as  $\mathcal{C}_0^{-1} = \mathcal{Y} + \eta\mathcal{Y}^2$  with  $\eta > 0$ . One can verify that this choice of prior covariance operator is symmetric, positive definite and trace-class as long as  $\gamma, \eta, \beta > 0$ . The noise covariance operator for the observations is chosen to be a multiple of the identity matrix, i.e.,  $\mathbf{\Gamma}_{\text{noise}} = \sigma^2 \mathbf{I}$ —in our examples we choose  $\sigma = 1$ . The parameters  $\gamma$ ,  $\beta$ , and  $\eta$  are chosen as  $\gamma = 10^{-3}$ ,  $\beta = 10^{-4}$  and  $\eta = 10^{-2}$ , and we have verified that this choice approximately satisfies the discrepancy principle. In the (discrete) numerical applications, we use  $\delta = 0$  in the measure  $\mu_\delta$  the trace estimator vectors  $z_i$  are sampled from (see section 4.2).

To enforce the constraint  $\mathbf{w} \in \mathcal{S}$ , i.e.,  $\sqrt{w_1^2 + w_2^2} = 1$  in this case, we parameterize the weight vector as  $(w_1, \pm\sqrt{1-w_1^2})$ . The parameter  $w_1$ , alone, controls the combination of both sources. Moreover, the weight vectors  $(w_1, -\sqrt{1-w_1^2})$  and  $(-w_1, \sqrt{1-w_1^2})$  lead to the same reconstruction, such that it suffices to consider the encoding weights  $(w_1, \sqrt{1-w_1^2})$  for  $w_1 \in [-1, 1]$ .

In figure 2, we plot the three objective functions (12)–(14) from section 4.2. For each  $w_1 \in [-1, 1]$ , the Gauss–Newton (12) and Laplace (13) formulations are evaluated at the MAP point,  $m_{\text{MAP}}(w_1)$ , corresponding to the encoding weight  $(w_1, \sqrt{1-w_1^2})$ ; in other words, the Hessian for these two criteria is evaluated at a medium parameter  $m_{\text{MAP}}(w_1)$  that varies with the weight  $w_1$ . For formulation (14), we choose  $m_0$  to be a constant value equal to the background medium, i.e.,  $m_0 \equiv 1$ . We observe that the result for the Gauss–Newton formulation (12) differs from the Laplace approximation (13). In addition, we clearly observe that each formulation contains local minima.

*Robustness of the Gauss–Newton formulation (12)* Since the computation of the MAP point  $m_{\text{MAP}}(w_1)$  is a computationally intensive task for large-scale problems, it

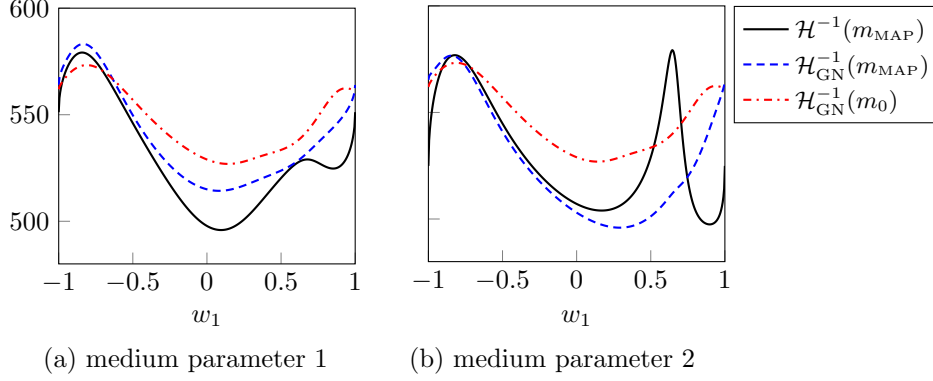


Figure 2: Plots of  $\text{tr}(\mathcal{H}^{-1})$  with  $\mathcal{H}^{-1}(m_{\text{MAP}}(w_1))$ ,  $\mathcal{H}_{\text{GN}}^{-1}(m_{\text{MAP}}(w_1))$  and  $\mathcal{H}_{\text{GN}}^{-1}(m_0)$  for both target media.  $m_0 \equiv 1$ , same as the background value for the medium parameter.

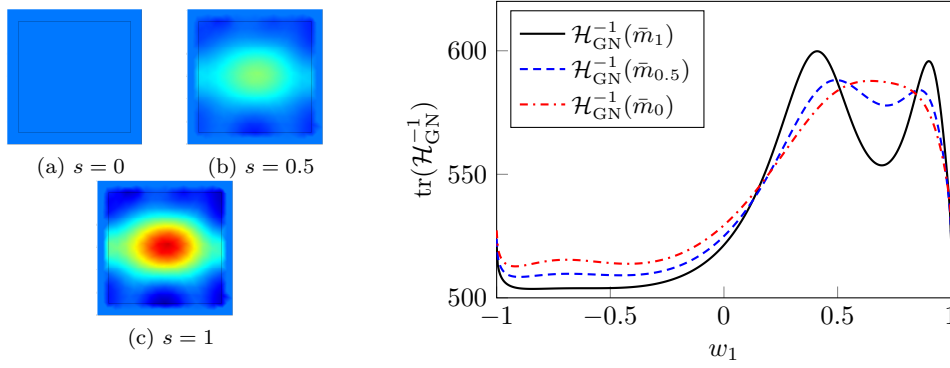


Figure 3: Plots of objective function  $\Phi_0$  (14) for weights  $w_1 \in [-1, 1]$  (right), at medium  $\bar{m}_s$ , with  $s = 0, 0.5, 1$  (left). Here  $m_0 \equiv 1$  (the background medium).

might be useful to solve the optimization (11) without having to recompute the exact MAP point for each iterate of the weights. The Laplace formulation (13) is based on the full Hessian which is guaranteed to be positive definite only in a neighbourhood of the MAP point. The Gauss–Newton approximation, however, is always positive definite and we observe numerically that it preserves relevant information about the objective function, even far away from the MAP point. In figure 3, we plot the objective function (12), for all values of  $w_1 \in [-1, 1]$ , for different (fixed) medium parameters  $\bar{m}_s$  ranging from the background medium,  $m_0 \equiv 1$ , to the MAP point  $m^\sharp$  computed using both sources independently (for medium parameter 2). The sources are located at the points  $(0, 0.1)$  and  $(0, 1.1)$ . That is, we define

$$\bar{m}_s = (1 - s)m_0 + s m^\sharp.$$

It appears that the medium parameter needs to include the main features of the target medium sufficiently accurately ( $s > 0.5$ ) to match the main features of the exact trace of the posterior covariance; this can be seen from the behavior of  $\text{tr}(\mathcal{H}_{\text{GN}}^{-1}(w_1, \bar{m}_s))$  in the interval  $w_1 \in [0.2, 1.0]$ .

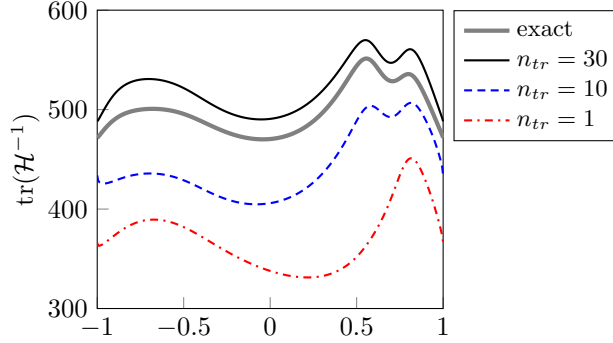


Figure 4: Plots of the objective function in (16) when the trace of the posterior covariance is computed exactly or with a trace estimator ( $n_{tr} = 1, 10, 30$ ). For each  $n_{tr}$ , we used a fixed realization of the trace estimator.

*The effect of trace estimation* When computing A-optimal encoding weights, one only needs the local minima of the trace to be well characterized. We show in figure 4 that trace estimation does indeed affect the shape of the objective function in the formulations of the A-optimal encoding weights (16). However, in our example, the objective function using a trace estimation preserves the local minima of the objective function using an exact trace when a sufficient number of trace estimator vectors are used.

### 6.2. A-optimal encoding weights in higher dimensional weight spaces

We now consider a problem with 10 sources (i.e.,  $N_s = 10$ ). Here, we focus

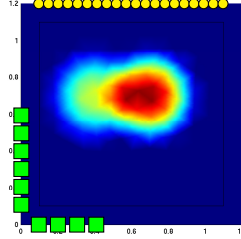


Figure 5: Target medium parameter and locations of the 10 sources (green squares), and receivers (yellow circles).

on qualitative properties of the A-optimal source encoding weights by performing statistical tests, in which we study how successful A-optimal encoding weights are in reducing posterior variance and relative medium misfit compared to encoding weights sampled from the uniform spherical distribution. We also compared with random weights sampled, then re-scaled, from the Rademacher distribution (see section 2). Since the results we obtained were not statistically different from the results presented in this section using random weights sampled from the uniform spherical distribution, we decided to omit these results. Throughout this section, the relative medium misfit is taken to be the relative  $L^2$ -error between the reconstruction of interest and the

reconstruction obtained using all 10 sources independently. The penalty parameter was empirically selected to be  $\lambda = 10^3$ .

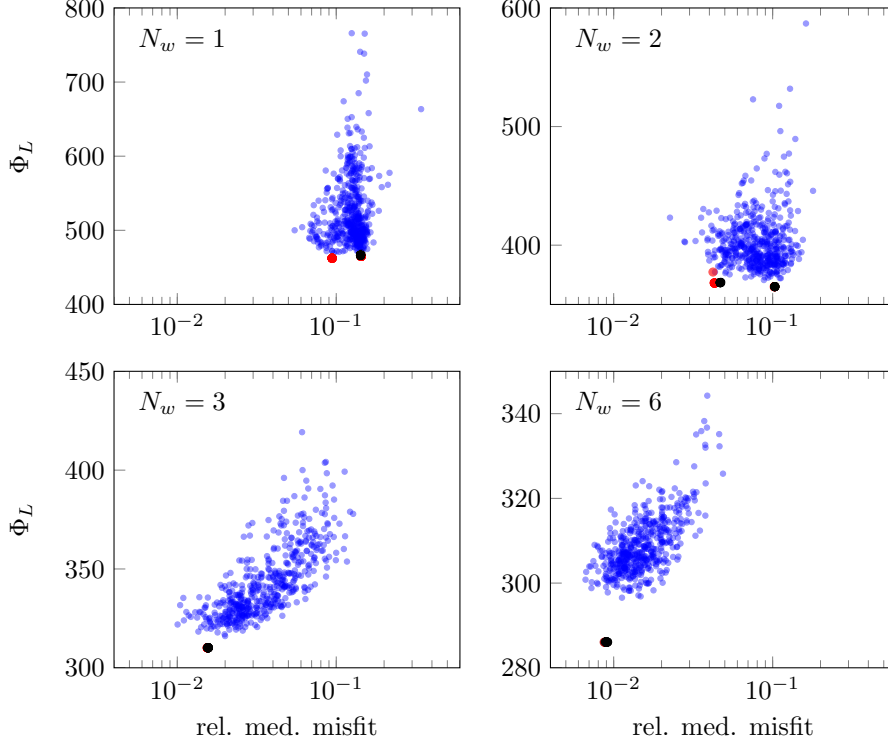


Figure 6: Plot of  $\Phi_L(\mathbf{w})$  (13) against relative medium misfit ( $N_s = 10$  and  $N_w = 1, 2, 3, 6$ ) for reconstructions using random encoding sources sampled from the uniform spherical distribution (blue) or A-optimal encoding weights computed with formulation (15) (black) and (16) (red). Target model 2 with source configuration as shown in figure 5. Sample size = 500,  $n_{tr} = 30$ .

We show the results in figure 6. Each plot shows, for different number of encoded sources ( $N_w = 1, 2, 3$  and  $6$ ), the objective function  $\Phi_L(\mathbf{w})$  defined in (13) against the relative medium misfit of the reconstruction, which is an indication for the quality of the reconstruction. Each reconstruction is indicated by a translucent dot; a darker shade indicates a higher concentration of reconstructions in that part of the plot. This shows the variation in the quality of the reconstruction. The blue dots correspond to reconstructions that use random encoding weights sampled from the uniform spherical distribution. The red dots indicate A-optimal encoding weights based on the Laplace formulation (16). The reconstructions marked with black dots use A-optimal encoding weights based on the Gauss-Newton formulation (15). In order to detect potential local minima, the A-optimal encoding weights are re-computed several times, starting from different initial conditions.

Notice that with one encoded source, A-optimal encoding weights do not provide a clear advantage over random weights. The overall distribution of random weights does not indicate a strong connection between the trace of the posterior covariance (13)

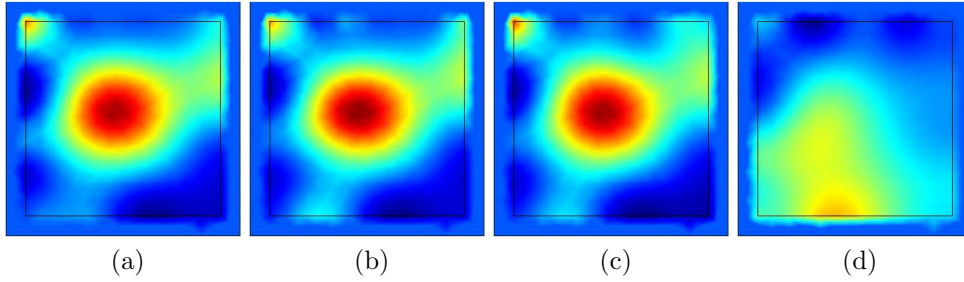


Figure 7: Four examples of reconstructions using different number of sources, with target parameter 2: (a) 10 independent sources; (b) 3 A-optimally encoded sources; (c) 3 randomly encoded sources; (d) 3 other randomly encoded sources.

and the relative medium misfit. On the other hand, the A-optimal encoding weights outperform the random weights (on average), when sufficiently many encoding weights are used (see in particular  $N_w = 2$  and 3 in figure 6). In that case, the random weights appear to indicate a linear correlation between our objective function and the relative medium misfit, which translates into the best reconstruction being also the one with smallest trace of the posterior covariance. Overall, these results suggest the existence of a threshold, in the number of encoding sources, above which optimal weights provide improvement in both variance and medium misfit over random encoding weights. Moreover, based on these results, there does not appear to be a clear advantage in using the Laplace approximation (16) over the Gauss–Newton approximation (15), provided sufficiently many encoded sources are used. In the last row of figure 6, optimal weights computed with both formulations provide similar results, although the actual values of the weights do not necessarily agree.

In addition, we provide a comparison of the reconstructions computed using all sources independently (figure 7a), using three A-optimally encoded sources (figure 7b), and two examples of reconstructions computed using three randomly encoded sources: one resulting in a good reconstruction (figure 7c), and one resulting in a poor reconstruction (figure 7d). There is virtually no difference between the reconstructions computed using all 10 sources and using three A-optimally encoded sources. On the other hand, using random encoding weights drawn from the same distribution may lead to good or poor reconstructions, as is shown in figures 7c, d. This is consistent with the results in figure 6 (bottom left), where the blue dots show large variations in terms of relative medium misfit.

*Variability of the A-optimal encoding weights* The A-optimal encoding weight formulation introduced in section 4 relies on a fixed realization of the trace estimator. Note that the A-optimal encoding weights are solutions to a highly nonlinear optimization problem that in general exhibits local minima. However, we show numerically that, provided sufficiently many encoding weights are chosen and a large enough number of trace estimators are used, the computation of the A-optimal encoding weights is stable with respect to trace estimation. In figure 8, we show 100 results obtained with Laplace A-optimal encoding weights (16), in the case of 3 encoded sources, with different numbers of trace estimators ( $n_{tr} = 4, 10, 30$ ). Each computation uses different realizations of the trace estimator, and different initial

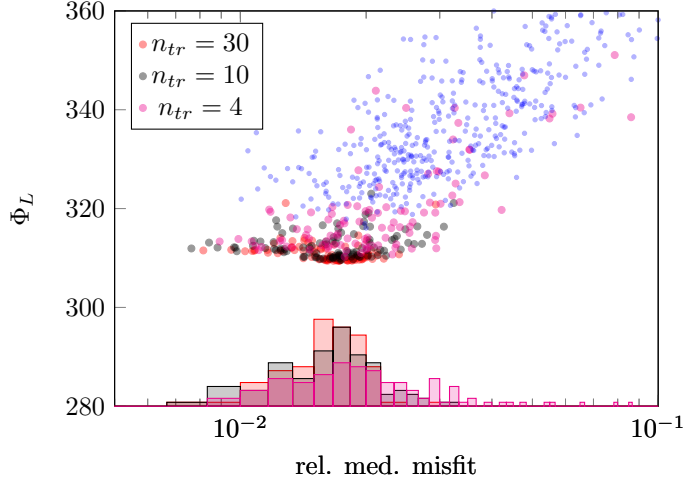


Figure 8: Variability of the A-optimal weights for different numbers of trace estimators,  $n_{tr} = 30$  (red), 10 (black) and 4 (magenta). A-optimal encoding weights are computed with formulation (16) ( $N_s = 10$  and  $N_w = 3$ ), using different realizations of the trace estimators and different initial guess of the weights. Sample size = 100.

guess of the weights.

We observe that with  $n_{tr} = 10$  and 30 the computations of the A-optimal encoding weights provide similar results. On the other hand, the use of 4 trace estimator vectors leads to a much wider range in the quality of the results, both in terms of relative medium misfit and trace of the posterior covariance.

### 6.3. Remarks on the Gauss–Newton formulation

Here, we discuss the justification for and advantages of using the Gauss–Newton formulation for finding A-optimal encoding weights. In many important situations, the Gauss–Newton formulation appears accurate enough to compute the A-optimal encoding weights. The Gauss–Newton approximation to the Hessian is most accurate when the data misfit residual is small at the solution of the inverse problem. This is the case, for instance, when the noise level in the observations is low. In our numerical experiments we observed that, provided sufficiently many encoded sources are used, the Gauss–Newton formulation represents a sufficiently accurate approximation to the Laplace formulation for the purpose of computing A-optimal encoding weights.

The Gauss–Newton formulation holds strong promises to reduce the computational cost of the A-optimal encoding weights. The data-misfit part of the Gauss–Newton Hessian is guaranteed to be positive semi-definite at any evaluation point, and hence the Gauss–Newton Hessian is positive definite. This allows two main improvements to the computations of the A-optimal weights. First, and as detailed in section 5.2.3, one can incorporate a low-rank approximation of the Gauss–Newton Hessian to reduce the computational cost. The magnitude of that reduction is problem-dependent, but will be most noticeable when large numbers of trace estimators are required.

Another advantage of the positive definiteness of the Gauss–Newton Hessian is that the objective function (12) of the Gauss–Newton formulation does not have to be evaluated in a small neighbourhood of the MAP point for the objective function to make sense. This could allow one, for instance, to solve the MAP point inexactly when the A-optimal objective function is far from its minimum, which would reduce the overall computational cost. In section 6.1, we studied how the objective function varies with the evaluation point  $\bar{m}_s$  (figure 3), and observed that the objective function tends to maintain similar local minima away from the MAP point.

Finally, we want to point out that in certain situations, the full Hessian may not be available, may be too complicated to derive, or too expensive to compute, rendering the Laplace formulation inadequate. This can be the case for inverse problems with highly nonlinear forward problems.

## 7. Conclusion

We have developed a method for the computation of A-optimal encoding weights aiming at large-scale non-linear inverse problems. As we show numerically, reconstructions obtained using A-optimal encoding weights not only minimize the average of the posterior variance, but consistently outperform random encoding weights in terms of the quality of the reconstructions. While in this work, we relied on quasi-Newton methods for solving the optimization problem for A-optimal encoding weights, we will explore the derivation and implementation of a Newton solver for this optimization problem in future work. We point out that, thanks to the optimize-then-discretize approach we adopted, the derivation of the analytical expression for the action of the Hessian in a direction is possible with little more effort than what was required to get the gradient.

We introduced two formulations for the computation of the A-optimal encoding weights, namely the Gauss–Newton formulation (15) and the Laplace formulation (16). Although the Gauss–Newton formulation represents an approximation to the Laplace formulation, it holds several advantageous features from computational point of view.

We note that computing A-optimal encoding weights can entail a significant computational effort. However, the method can be attractive for real-time monitoring applications where one needs to solve an inverse problem repeatedly over time. In this case, one first computes the A-optimal encoding weights offline, and then can use those weights to solve the inverse problem repeatedly at a fraction of the original cost. An example for such an application is the monitoring of an oil reservoir, where seismic or electro-magnetic inverse problems are solved repeatedly to characterize the evolution of the reservoir properties over time.

## References

- [1] Fredi Tröltzsch. *Optimal Control of Partial Differential Equations: Theory, Methods and Applications*, volume 112 of *Graduate Studies in Mathematics*. American Mathematical Society, 2010.
- [2] Jerome R. Krebs, John E. Anderson, David Hinkley, Ramesh Neelamani, Sunwoong Lee, Anatoly Baumstein, and Martin-Daniel Lacasse. Fast full-wavefield seismic inversion using encoded sources. *Geophysics*, 74(6):WCC177–WCC188, 2009.
- [3] Partha S Routh, Sunwoong Lee, Ramesh Neelamani, Jerome R Krebs, Spyridon Lazaratos, and Carey Marcinkovich. Simultaneous source encoding and source separation as a practical solution for full wavefield inversion, September 9 2011. US Patent App. 13/229,252.



- [4] Eldad Haber, Matthias Chung, and Felix J Herrmann. An effective method for parameter estimation with PDE constraints with multiple right hand sides. *SIAM Journal on Optimization*, 22 (3):739–757, 2012.
- [5] Ellen B. Le, Aaron Myers, and Tan Bui-Thanh. A Randomized Misfit Approach for Data Reduction in Large-Scale Inverse Problems. *ArXiv e-prints*, March 2016.
- [6] Eldad Haber and Matthias Chung. Simultaneous source for non-uniform data variance and missing data. *CoRR*, abs/1404.5254, 2014.
- [7] William W. Symes. Source synthesis for waveform inversion. Technical report, Rice University, CAM report TR10-12, 2010.
- [8] Eldad Haber, Kees van den Doel, and Lior Horesh. Optimal design of simultaneous source encoding. *Inverse Problems in Science and Engineering*, pages 1–18, 2014.
- [9] Alen Alexanderian, Noemi Petra, Georg Stadler, and Omar Ghattas. A-optimal design of experiments for infinite-dimensional Bayesian linear inverse problems with regularized  $\ell_0$ -sparsification. *SIAM Journal on Scientific Computing*, 36(5):A2122–A2148, 2014.
- [10] Alen Alexanderian, Noemi Petra, Georg Stadler, and Omar Ghattas. A fast and scalable method for A-optimal design of experiments for infinite-dimensional Bayesian nonlinear inverse problems. *SIAM Journal on Scientific Computing*, 38(1):A243–A272, 2016.
- [11] Eldad Haber, Lior Horesh, and Luis Tenorio. Numerical methods for experimental design of large-scale linear ill-posed inverse problems. *Inverse Problems*, 24(055012):125–137, 2008.
- [12] Eldad Haber, Lior Horesh, and Luis Tenorio. Numerical methods for the design of large-scale nonlinear discrete ill-posed inverse problems. *Inverse Problems*, 26(2):025002, 2010.
- [13] Dariusz Uciński. *Optimal measurement methods for distributed parameter system identification*. CRC Press, Boca Raton, 2005.
- [14] Loyd N. Trefethen and David Bau. *Numerical Linear Algebra*. SIAM, Philadelphia, 1997.
- [15] Haim Avron and Sivan Toledo. Randomized algorithms for estimating the trace of an implicit symmetric positive semi-definite matrix. *Journal of the ACM (JACM)*, 58(2):17, April 2011.
- [16] Michael F. Hutchinson. A stochastic estimator of the trace of the influence matrix for Laplacian smoothing splines. *Communications in Statistics-Simulation and Computation*, 19(2):433–450, 1990.
- [17] Theodore W. Anderson and Michael A. Stephens. Tests for randomness of directions against equatorial and bimodal alternatives. *Biometrika*, 59(3):613–621, 1972.
- [18] Alexander Shapiro, Darinka Dentcheva, and Andrej Ruszczyński. *Lectures on Stochastic Programming: Modeling and Theory*. Society for Industrial and Applied Mathematics, 2009.
- [19] Andrew M. Stuart. Inverse problems: A Bayesian perspective. *Acta Numerica*, 19:451–559, 2010.
- [20] Masoumeh Dashti and Andrew M. Stuart. The Bayesian approach to inverse problems. In Roger Ghanem, David Higdon, and Houman Owhadi, editors, *Handbook of Uncertainty Quantification*. Springer, 2015.
- [21] Tan Bui-Thanh, Omar Ghattas, James Martin, and Georg Stadler. A computational framework for infinite-dimensional Bayesian inverse problems Part I: The linearized case, with application to global seismic inversion. *SIAM Journal on Scientific Computing*, 35(6):A2494–A2523, 2013.
- [22] 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.
- [23] Friedrich Pukelsheim. *Optimal Design of Experiments*. John Wiley & Sons, New-York, 1993.
- [24] Anthony C. Atkinson and Alexander N. Donev. *Optimum Experimental Designs*. Oxford, 1992.
- [25] Kathryn Chaloner and Isabella Verdinelli. Bayesian experimental design: A review. *Statistical Science*, 10(3):273–304, 1995.
- [26] Alfio Borzi and Volker Schulz. *Computational Optimization of Systems Governed by Partial Differential Equations*. SIAM, 2012.
- [27] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer Verlag, Berlin, Heidelberg, New York, second edition, 2006.
- [28] Pearl H. Flath, Lucas C. Wilcox, Volkan Akçelik, Judy Hill, Bart van Bloemen Waanders, and Omar Ghattas. Fast algorithms for Bayesian uncertainty quantification in large-scale linear inverse problems based on low-rank partial Hessian approximations. *SIAM Journal on Scientific Computing*, 33(1):407–432, 2011.
- [29] Peter Deufhard. *Newton methods for nonlinear problems*, volume 35 of *Springer Series in Computational Mathematics*. Springer-Verlag, Berlin, 2004.

**Appendix A. Gradient of the optimization formulation (23)**

We detail the derivation of the gradient of the Laplace formulation of the A-optimal weights in the case of the Helmholtz inverse problem, as defined in (23). In that formulation, we enforce the PDE constraints weakly using Lagrange multipliers. Therefore, we need to introduce adjoint variables that are indicated with a star superscript, e.g.,  $m^*$  is the adjoint variable for  $m$ . Following the formal Lagrangian approach [1], we define the Lagrangian  $\mathcal{L}$ ,

$$\begin{aligned}
\mathcal{L}(\mathbf{w}, m, \{u_i\}, \{p_i\}, \{v_{i,k}\}, \{q_{i,k}\}, \{y_k\}, m^*, \{u_i^*\}, \{p_i^*\}, \{v_{i,k}^*\}, \{q_{i,k}^*\}, \{y_k^*\}) = \\
\frac{1}{n_{tr}} \sum_{k=1}^{n_{tr}} \langle y_k, z_k \rangle + \\
\frac{1}{n_{tr} N_w} \sum_{k=1}^{n_{tr}} \sum_{i=1}^{N_w} [\langle \nabla v_{i,k}, \nabla v_{i,k}^* \rangle - \kappa^2 \langle m v_{i,k}, v_{i,k}^* \rangle - \kappa^2 \langle u_i y_k, v_{i,k}^* \rangle] \\
+ \frac{1}{n_{tr} N_w} \sum_{k=1}^{n_{tr}} \sum_{i=1}^{N_w} \left[ \langle \nabla q_{i,k}^*, \nabla q_{i,k} \rangle - \kappa^2 \langle q_{i,k}^*, m q_{i,k} \rangle - \kappa^2 \langle q_{i,k}^*, p_i y_k \rangle + \langle B q_{i,k}^*, B v_{i,k} \rangle_{\mathbf{\Gamma}_{noise}^{-1}} \right] \\
+ \frac{1}{n_{tr}} \sum_{k=1}^{n_{tr}} \left[ \langle y_k, y_k^* \rangle_{\mathcal{E}} - \frac{1}{N_w} \sum_{i=1}^{N_w} \kappa^2 (\langle v_{i,k} p_i, y_k^* \rangle + \langle u_i q_{i,k}, y_k^* \rangle) - \langle z_k, y_k^* \rangle \right] + \\
\frac{1}{N_w} \sum_{i=1}^{N_w} [\langle \nabla u_i, \nabla u_i^* \rangle - \kappa^2 \langle m u_i, u_i^* \rangle - \langle f(\mathbf{w}^i), u_i^* \rangle] \\
+ \frac{1}{N_w} \sum_{i=1}^{N_w} \left[ \langle \nabla p_i^*, \nabla p_i \rangle - \kappa^2 \langle p_i^*, m p_i \rangle + \langle B p_i^*, B u_i - \mathbf{d}(\mathbf{w}^i) \rangle_{\mathbf{\Gamma}_{noise}^{-1}} \right] \\
+ \langle m - m_0, m^* \rangle_{\mathcal{E}} - \frac{1}{N_w} \sum_{i=1}^{N_w} \kappa^2 \langle u_i p_i, m^* \rangle. \tag{A.1}
\end{aligned}$$

The gradient is then given by  $\delta_{\mathbf{w}} \mathcal{L} = [\delta_{\mathbf{w}^1} \mathcal{L}, \delta_{\mathbf{w}^2} \mathcal{L}, \dots, \delta_{\mathbf{w}^{N_w}} \mathcal{L}]^T$ , where for any  $i = 1, \dots, N_w$ ,

$$\delta_{\mathbf{w}^i} \mathcal{L} = -\frac{1}{N_w} \begin{bmatrix} \langle f_1, u_i^* \rangle + \langle B p_i^*, \mathbf{d}_1 \rangle_{\mathbf{\Gamma}_{noise}^{-1}} \\ \langle f_2, u_i^* \rangle + \langle B p_i^*, \mathbf{d}_2 \rangle_{\mathbf{\Gamma}_{noise}^{-1}} \\ \vdots \\ \langle f_{N_s}, u_i^* \rangle + \langle B p_i^*, \mathbf{d}_{N_s} \rangle_{\mathbf{\Gamma}_{noise}^{-1}} \end{bmatrix}.$$

Before we specify the steps that lead to the evaluation of the variables  $u_i^*$  and  $p_i^*$ , we identify some important symmetries between the state variables and their adjoints. Indeed, for each  $k = 1, \dots, n_{tr}$ , the variables  $(y_k \{v_{i,k}\}_i, \{q_{i,k}\}_i)$  solve a Hessian system similar to (22), and the corresponding adjoint variables  $(y_k^* \{v_{i,k}^*\}_i, \{q_{i,k}^*\}_i)$  solve the system of equations given (formally) by  $\delta_{v_{i,k}} \mathcal{L} = \delta_{q_{i,k}} \mathcal{L} = \delta_{y_k} \mathcal{L} = 0$ . While the former system of equations solve  $\mathcal{H} y_k = z_k$ , the latter solves  $\mathcal{H} y_k^* = -z_k$ . This leads to the symmetry relations

$$y_k = -y_k^*, \quad v_{i,k} = -v_{i,k}^*, \quad \text{and} \quad q_{i,k} = -q_{i,k}^*, \tag{A.2}$$

for any  $i = 1, \dots, N_w$  and  $k = 1, \dots, n_{tr}$ .

For any  $i = 1, \dots, N_w$ , the variable  $u_i^*$  (resp.  $p_i^*$ ) solves the equation  $\delta_{u_i} \mathcal{L} = 0$  (resp.  $\delta_{p_i} \mathcal{L} = 0$ ). That is, for any  $\tilde{u} \in H^1(\mathcal{D})$ ,  $u_i^*$  solves

$$\begin{aligned} & \langle \nabla u_i^*, \nabla \tilde{u} \rangle - \kappa^2 \langle m u_i^*, \tilde{u} \rangle \\ & - \kappa^2 \langle p_i m^*, \tilde{u} \rangle + \langle B p_i^*, B \tilde{u} \rangle_{\Gamma_{\text{noise}}^{-1}} - \kappa^2 \frac{1}{n_{tr}} \sum_{k=1}^{n_{tr}} [\langle y_k v_{i,k}^*, \tilde{u} \rangle + \langle q_{i,k} y_k^*, \tilde{u} \rangle] = 0. \end{aligned}$$

On the other hand, for any  $\tilde{p} \in H^1(\mathcal{D})$ ,  $p_i^*$  solves

$$\langle \nabla p_i^*, \nabla \tilde{p} \rangle - \kappa^2 \langle p_i^*, m \tilde{p} \rangle - \kappa^2 \langle u_i m^*, \tilde{p} \rangle - \kappa^2 \frac{1}{n_{tr}} \sum_{k=1}^{n_{tr}} [\langle q_{i,k}^* y_k, \tilde{p} \rangle + \langle v_{i,k} y_k^*, \tilde{p} \rangle] = 0.$$

Using (A.2), this reduces, for any  $i = 1, \dots, N_w$ , to the system of equations

$$\begin{aligned} & \langle \nabla u_i^*, \nabla \tilde{u} \rangle - \kappa^2 \langle m u_i^*, \tilde{u} \rangle - \kappa^2 \langle p_i m^*, \tilde{u} \rangle + \langle B p_i^*, B \tilde{u} \rangle_{\Gamma_{\text{noise}}^{-1}} + \frac{2}{n_{tr}} \sum_{k=1}^{n_{tr}} \kappa^2 \langle y_k q_{i,k}, \tilde{u} \rangle = 0, \\ & \langle \nabla p_i^*, \nabla \tilde{p} \rangle - \kappa^2 \langle m p_i^*, \tilde{p} \rangle - \kappa^2 \langle u_i m^*, \tilde{p} \rangle + \frac{2}{n_{tr}} \sum_{k=1}^{n_{tr}} \kappa^2 \langle v_{i,k} y_k, \tilde{p} \rangle = 0. \end{aligned} \tag{A.3}$$

Therefore, computation of the  $u_i^*$ 's and  $p_i^*$ 's requires knowledge of the quantities  $\{u_i\}$ ,  $\{p_i\}$ ,  $m^*$ ,  $\{v_{i,k}\}$ ,  $\{q_{i,k}\}$  and  $\{y_k\}$ . Variables  $\{u_i\}$ ,  $\{p_i\}$ ,  $\{v_{i,k}\}$ ,  $\{q_{i,k}\}$ , and  $\{y_k\}$  are all evaluated during the computation of the objective functional  $1/n_{tr} \sum_{k=1}^{n_{tr}} \langle y_k, z_k \rangle$ , such that the only remaining unknown quantity is  $m^*$ . That variable is solution to the equation  $\delta_m \mathcal{L} = 0$ , that is, for any  $\tilde{m} \in \mathcal{E}$ ,  $m^*$  solves

$$\begin{aligned} & \frac{1}{n_{tr} N_w} \sum_{k=1}^{n_{tr}} \sum_{i=1}^{N_w} [-\kappa^2 \langle \tilde{m} v_{i,k}, v_{i,k}^* \rangle - \kappa^2 \langle q_{i,k}^*, \tilde{m} q_{i,k} \rangle] \\ & + \frac{1}{N_w} \sum_{i=1}^{N_w} [-\kappa^2 \langle \tilde{m} u_i, u_i^* \rangle - \kappa^2 \langle p_i^*, \tilde{m} p_i \rangle] + \langle \tilde{m}, m^* \rangle_{\mathcal{E}} = 0. \end{aligned}$$

Using (A.2), we simplify this equation to obtain

$$\frac{2}{n_{tr} N_w} \sum_{k=1}^{n_{tr}} \sum_{i=1}^{N_w} \kappa^2 \langle v_{i,k} q_{i,k}, \tilde{m} \rangle - \frac{1}{N_w} \sum_{i=1}^{N_w} \kappa^2 [\langle u_i u_i^*, \tilde{m} \rangle + \langle p_i^* p_i, \tilde{m} \rangle] + \langle m^*, \tilde{m} \rangle_{\mathcal{E}} = 0.$$

This equation can be grouped with the system of equations (A.3) to obtain the larger system

$$\begin{aligned} & \langle \nabla p_i^*, \nabla \tilde{p} \rangle - \kappa^2 \langle m p_i^*, \tilde{p} \rangle - \kappa^2 \langle u_i m^*, \tilde{p} \rangle = -\frac{2}{n_{tr}} \sum_{k=1}^{n_{tr}} \kappa^2 \langle v_{i,k} y_k, \tilde{p} \rangle \\ & \langle \nabla u_i^*, \nabla \tilde{u} \rangle - \kappa^2 \langle m u_i^*, \tilde{u} \rangle - \kappa^2 \langle p_i m^*, \tilde{u} \rangle + \langle B p_i^*, B \tilde{u} \rangle_{\Gamma_{\text{noise}}^{-1}} = -\frac{2}{n_{tr}} \sum_{k=1}^{n_{tr}} \kappa^2 \langle y_k q_{i,k}, \tilde{u} \rangle \\ & \langle m^*, \tilde{m} \rangle_{\mathcal{E}} - \frac{1}{N_w} \sum_{i=1}^{N_w} \kappa^2 [\langle u_i u_i^*, \tilde{m} \rangle + \langle p_i^* p_i, \tilde{m} \rangle] = -\frac{2}{n_{tr} N_w} \sum_{k=1}^{n_{tr}} \sum_{i=1}^{N_w} \kappa^2 \langle v_{i,k} q_{i,k}, \tilde{m} \rangle. \end{aligned}$$

This system of equations should be compared to the system of equations for the Hessian (22). From this, it should be clear that the computation of  $m^*$  corresponds to the solution of another Hessian system with a right-hand side depending on the state and adjoint variables,  $\{u_i\}$  and  $\{p_i\}$ , the incremental state and adjoint variables,  $\{v_{i,k}\}$  and  $\{q_{i,k}\}$ , the medium parameter  $m$ , and the  $\{y_k\}$ . We denote this right-hand side as  $\mathcal{F}$ . In strong form,  $m^*$  thus solves

$$\mathcal{H}(m_{\text{MAP}})m^* = \mathcal{F}(\{u_i\}, \{p_i\}, \{v_{i,k}\}, \{q_{i,k}\}, m, \{y_k\}).$$