

A greedy sensor selection algorithm for hyperparameterized linear Bayesian inverse problems with correlated noise models

Nicole Aretz^a, Peng Chen^b, Denise Degen^c, Karen Veroy^{d,*}

^a Oden Institute for Computational Engineering and Sciences, University of Texas at Austin, 201 E 24th St, Austin, TX 78712, USA

^b School of Computational Science and Engineering, Georgia Institute of Technology, 756 W Peachtree St NW, Atlanta, GA 30308, USA

^c Computational Geoscience, Geothermics, and Reservoir Geophysics, RWTH Aachen University, Mathieustr. 30, 52074 Aachen, Germany

^d Center for Analysis, Scientific Computing and Applications, Department of Mathematics and Computer Science, Eindhoven University of Technology, 5612 AZ Eindhoven, the Netherlands

ARTICLE INFO

Keywords:

Optimal sensor placement
Bayesian inverse problems
Correlated noise
Model order reduction
Greedy algorithm
Orthogonal matching pursuit

ABSTRACT

We consider optimal sensor placement for a family of linear Bayesian inverse problems characterized by a deterministic hyper-parameter. The hyper-parameter describes distinct configurations in which measurements can be taken of the observed physical system. To optimally reduce the uncertainty in the system's model with a single set of sensors, the initial sensor placement needs to account for the non-linear state changes of all admissible configurations. We address this requirement through an observability coefficient which links the posteriors' uncertainties directly to the choice of sensors. We propose a greedy sensor selection algorithm to iteratively improve the observability coefficient for all configurations through orthogonal matching pursuit. The algorithm allows explicitly correlated noise models even for large sets of candidate sensors, and remains computationally efficient for high-dimensional forward models through model order reduction. We demonstrate our approach on a large-scale geophysical model of the Perth Basin, and provide numerical studies regarding optimality and scalability with regard to classic optimal experimental design utility functions.

1. Introduction

In the Bayesian approach to inverse problems (cf. [1]), the uncertainty in a parameter is described via a probability distribution. With Bayes' Theorem, the prior belief in a parameter is updated when new information is revealed such that the posterior distribution describes the parameter with improved certainty. Bayes' posterior is optimal in the sense that it is the unique minimizer of the sum of the relative entropy between the posterior and the prior, and the mean squared error between the model prediction and the experimental data ([2,3]). The noise model drives, along with the measurements, how the posterior's uncertainty is reduced in comparison to the prior. A critical aspect – especially for expensive experimental data – is how to select the measurements to improve the posterior's credibility best. For example, when simulating the subsurface heat distribution for geothermal applications, unknown parameters (e.g., the geothermal heat flux, see Section 5) need to be inferred from temperature measurements. The measurements are

* Corresponding author.

E-mail address: k.p.veroy@tue.nl (K. Veroy).

taken in boreholes, which can cost several million dollars to drill, so it is essential to plan their location carefully. To provide some perspective, the developments costs of a geothermal project (e.g., drilling, stimulation, and tests) take up 50–70% of the total budget ([4]). The selection of adequate sensors meeting individual applications' needs is, therefore, a big goal of the optimal experimental design (OED) research field and its surrounding community. We refer to the literature (e.g., [5–7]) for introductions.

In this paper, we consider inverse problem settings in which a deterministic hyper-parameter describes anticipated system configurations such as material properties or loading conditions. Each configuration changes the model non-linearly, so we obtain a *family* of possible posterior distributions for any measurement data. Supposing data can only be obtained with a single set of sensors regardless of the system's configuration, the OED task becomes to reduce the uncertainty in each posterior uniformly over all hyper-parameters. This task is challenging 1) for high-dimensional models since each configuration requires its own computationally expensive model solve, and 2) for correlated noise models since the non-nested structure of the inverse noise covariance matrix can cause discontinuities in relaxed, weighting-based approaches [8]. By building upon [9], this paper addresses both challenges and proposes a novel sensor selection algorithm that remains efficient even for large sets of admissible measurements. For instance, in Section 5.4 we apply our algorithm to a geophysical problem with 132,651 degrees of freedom in the state variable and 11,045 available sensor positions with correlated observations, both of which are high-dimensional.

The concept of the hyper-parameter is similar to so-called *nuisance* parameters in the literature. Nuisance parameters are a secondary source of uncertainty, causing additional variability in the measurements, while not being of primary interest for the inversion. Neglecting this source of uncertainty in the inverse problem can cause serious overconfidence in the inferred parameter while inverting for all uncertain parameters together increases the computational burden (see [10] for a comparison). In practice, the nuisance parameters are therefore often marginalized out with the Bayesian approximation error approach (e.g., [11,12]). Marginalization has also been adopted for OED over models with nuisance parameters, see [13] for A-optimal experimental design (A-OED), [14,15] for the expected information gain (EIG). For E-optimal experimental design (E-OED), [16] keeps both uncertain parameters in the inverse problem, but poses its OED formulation over a submatrix of the Fisher information matrix. Albeit the hyper-parameter we consider here can be interpreted as a form of model uncertainty, it differs from nuisance parameters through its aleatoric nature: its uncertainty cannot be reduced in the inverse problem, and we are therefore treating it as a separate, deterministic parameter. With this interpretation, our baseline setting concurs with [17] about A-OED under irreducible model uncertainty. However, in [17] the model uncertainty is integrated out by taking the expectation over the utility function, thereby obtaining a risk-neutral design that is favorable for most model realizations. In contrast, we are optimizing the minimum of the utility function to guarantee that the design remains informative in the worst-case scenario, which is often the requirement for risk-averse applications.

The main contributions are as follows: First, we identify an observability coefficient as a link between the sensor choice and the maximum eigenvalue of each posterior distribution. We provide an analysis of its sensitivity to model and parameter approximations. Second, we decompose the noise covariance matrix for any observation operator to allow fast computation of the observability coefficient's increase under expansion with additional sensors. The decomposition allows us to treat correlated noise covariance matrices efficiently when comparing the benefits of including additional sensors. Third, we propose a sensor selection algorithm that iteratively constructs an observation operator from a large set of sensors to increase the observability coefficient over all hyper-parameters. The algorithm is applicable to correlated noise models, and requires, through the efficient use of model order reduction (MOR) techniques, only a *single* full-order model evaluation per selected sensor, which achieves considerable efficiency.

The analysis and algorithm presented in this work significantly extend our initial ideas presented in [9] in which we seek to generalize the 3D-VAR stability results from [18] to the probabilistic Bayesian setting: This work additionally features 1) an analysis of the observability coefficient regarding model approximations, 2) explicit computational details for treating correlated noise models, and 3) a comprehensive discussion of the individual steps in the sensor selection algorithm. Moreover, the proposed method is tested using a large-scale geophysical model of the Perth Basin. Our proposed algorithm is directly related to the orthogonal matching pursuit (OMP) algorithm [19,20] for the parameterized-background data-weak (PBDW) method and the empirical interpolation method (EIM) ([21,22]). Closely related OED methods for linear Bayesian inverse problems over partial differential equations (PDEs) include [23–26,13,27], mostly for A- and D-OED and uncorrelated noise. In recent years, these methods have also been extended to non-linear Bayesian inverse problems, e.g., [28–32], while an advance to correlated noise has been made in [8]. In particular, [31,32] use similar algorithmic approaches to this work by applying a greedy algorithm to maximize the expected information gain. Common strategies for dealing with the high dimensions imposed by the PDE model use the framework in [33] for discretization, combined with parameter reduction methods (e.g., [34–40]) and MOR methods for uncertainty quantification (UQ) problems (e.g., [41–45]).

This paper is structured as follows: In Section 2 we introduce the hyper-parameterized inverse problem setting, including all assumptions for the prior distribution, the noise model, and the forward model. In Section 3, we then establish and analyze the connection between the observability coefficient and the posterior uncertainty. We finally propose our sensor selection algorithm in Section 4 which exploits the presented analysis to choose sensors that improve the observability coefficient even in a hyper-parameterized setting. In Section 5, we demonstrate the applicability and scalability of our approach on a geophysical model with high-dimensional state space before concluding in Section 6.

2. Problem setting

Let \mathcal{U} be a Hilbert space with inner product $\langle \cdot, \cdot \rangle_{\mathcal{U}}$ and induced norm $\|u\|_{\mathcal{U}}^2 := \langle u, u \rangle_{\mathcal{U}}$. We consider the problem of identifying unknown states $u = u_{\theta} \in \mathcal{U}$ of a single physical system under changeable configurations θ from noisy measurements

$$\mathbf{d}(\theta) \approx [\ell_1(u_\theta), \dots, \ell_K(u_\theta)]^T \in \mathbb{R}^K.$$

The measurements are obtained by a set of K unique *sensors* (or *experiments*) $\ell_1, \dots, \ell_K : \mathcal{U} \rightarrow \mathbb{R}$. Our goal is to choose these sensors from a large *sensor library* \mathcal{L} of options in a way that optimizes how much information is gained from their measurements for any configurations θ . In the following we specify our assumptions and provide the mathematical background to our setup.

Hyper-parameterized forward model

We consider the unknown state u to be uniquely characterized by two sources of information:

- an unknown parameter $\mathbf{m} \in \mathbb{R}^M$ describing uncertainties in the governing physical laws, and
- a hyper-parameter $\theta \in \mathcal{P} \subset \mathbb{R}^p$ describing dependencies on configurations under which the system may be observed (such as material properties or loading conditions) where \mathcal{P} is a given compact set enclosing all possible configurations. We interchangeably call θ *hyper-parameter* or *configuration* to either stress its role in the mathematical model or physical interpretation. We assume that potential variability in θ cannot be reduced in an inverse problem in that θ is either known and fixed in any inverse problems (e.g., when optimizing the geometry of mechanical parts in an outer loop to reduce the region of failure), or that θ describes irreducible model uncertainty as considered in [17]. This property distinguishes the hyper-parameter θ from *nuisance* parameters, which are uncertain parameters whose uncertainty can indeed be reduced in the inverse problem but which are not of primary interest. We refer to [46], section 1.3.2 for a more detailed distinction.

For any given $\mathbf{m} \in \mathbb{R}^M$ and $\theta \in \mathcal{P}$, we let $u_\theta(\mathbf{m}) \in \mathcal{U}$ be the solution of an abstract model equation $\mathcal{M}_\theta(u_\theta(\mathbf{m}); \mathbf{m}) = 0$ and assume that the map $\mathbf{m} \rightarrow u_\theta(\mathbf{m})$ is well-defined, linear, and uniformly continuous in \mathbf{m} , i.e.

$$\exists \bar{\eta} > 0 : \quad \bar{\eta}(\theta) := \sup_{\mathbf{m} \in \mathbb{R}^M} \frac{\|u_\theta(\mathbf{m})\|_{\mathcal{U}}}{\|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}}} < \bar{\eta} \quad \forall \theta \in \mathcal{P}. \quad (1)$$

Remark 1. Although we assumed that \mathbf{m} lies in the Euclidean space \mathbb{R}^M , any other finite-dimensional linear space can be considered via an affine transformation onto an appropriate basis (see [23,47]) while infinite-dimensional parameter spaces may be considered after appropriate discretization (cf. [33]). We note, however, that the numerical restrictions of our sensor placement algorithm in Section 3.3 require M to be small compared to the total number of sensors to be chosen. If this is not the case, we suggest to restrict the parameter space first onto a subspace spanned by the most uncertain parameter directions, or, for a goal-oriented experimental design, onto an active subspace (cf. [37,48,49,36]).

Remark 2. By keeping the model equation general, we stress the applicability of our approach to a wide range of problems. For instance, time-dependent states can be treated by choosing \mathcal{U} as a Bochner space or its discretization (cf. [50]). We also do not formally restrict the dimension of \mathcal{U} , though any implementation relies on the ability to compute $u_\theta(\mathbf{m})$ with sufficient accuracy. To this end, we note that the analysis in Section 3.2 can be applied to determine how discretization errors affect the observability criterion in the sensor selection.

Following a probabilistic approach to inverse problems, we express the initial uncertainty in $\mathbf{m} = \mathbf{m}(\theta)$ of any $u = u_\theta(\mathbf{m})$ in configuration θ through a random variable \mathbf{m} with Gaussian prior $\mu_{\text{pr}} = \mathcal{N}(\mathbf{m}_{\text{pr}}, \Sigma_{\text{pr}})$, where $\mathbf{m}_{\text{pr}} \in \mathbb{R}^M$ is the prior mean and $\Sigma_{\text{pr}} \in \mathbb{R}^{M \times M}$ is a symmetric positive definite (s.p.d.) covariance matrix. The latter defines the inner product $\langle \cdot, \cdot \rangle_{\Sigma_{\text{pr}}^{-1}}$ and its induced norm $\|\cdot\|_{\Sigma_{\text{pr}}^{-1}}$ through

$$\langle \mathbf{m}, \mathbf{v} \rangle_{\Sigma_{\text{pr}}^{-1}} := \mathbf{m}^T \Sigma_{\text{pr}}^{-1} \mathbf{v}, \quad \|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}}^2 := \langle \mathbf{m}, \mathbf{m} \rangle_{\Sigma_{\text{pr}}^{-1}}, \quad \forall \mathbf{m}, \mathbf{v} \in \mathbb{R}^M. \quad (2)$$

With these definitions, the probability density function (pdf) for μ_{pr} is

$$\pi_{\text{pr}}(\mathbf{m}) = \frac{1}{\sqrt{(2\pi)^M \det \Sigma_{\text{pr}}}} \exp\left(-\frac{1}{2} \|\mathbf{m} - \mathbf{m}_{\text{pr}}\|_{\Sigma_{\text{pr}}^{-1}}^2\right).$$

For simplicity, we assume $\{\mathbf{m}(\theta)\}_{\theta \in \mathcal{P}}$ to be independent realizations of \mathbf{m} such that we may consider the same prior for all θ without accounting for a possible history of measurements at different configurations.

Sensor library and noise model

For taking measurements of the unknown states $\{u(\theta)\}_\theta$, we call any linear functional $\ell \in \mathcal{U}'$ a *sensor*, and its application to a state $u \in \mathcal{U}$ its *measurement* $\ell(u) \in \mathbb{R}$. This implies in particular that any measurement $\ell(u)$ is linear in the state variable and bounded proportionally to the norm of the measured state $|\ell(x)| \leq \|\ell\|_{\mathcal{U}'} \|x\|_{\mathcal{U}}$. For the OED problem, we call the set $\mathcal{L} \subset \mathcal{U}'$ of admissible sensor choices our *sensor library*.

Example 1. Let $x : \Omega \rightarrow \mathbb{R}$ denote the temperature of a 2D domain $\Omega \subset \mathbb{R}^2$ modeling the surface of a work piece. A sensor taking a local temperature measurement at point $P = (P_1, P_2) \in \Omega$ can then be modeled as $\ell(x) = \frac{1}{\pi r^2} \int_{B_r(P)} x(s) ds$ where $B_r(P)$ is the ball

around P with probe radius $r > 0$. The sensor library might then be the corresponding set of sensors for all points P at which a probe may be placed on the physical asset.

We model noisy experimental measurements $d_\ell \in \mathbb{R}$ of the actual physical state u as $d_\ell = \ell(u) + \varepsilon_\ell$ where $\varepsilon_\ell \sim \mathcal{N}(0, \mathbf{cov}(\varepsilon_\ell, \varepsilon_\ell))$ is a Gaussian random variable. We permit noise in different sensor measurements to be correlated with a known covariance function \mathbf{cov} . In a slight overload of notation, we write $\mathbf{cov} : \mathcal{L} \times \mathcal{L} \rightarrow \mathbb{R}$, $\mathbf{cov}(\ell_i, \ell_j) := \mathbf{cov}(\varepsilon_{\ell_i}, \varepsilon_{\ell_j})$ as a symmetric bilinear form over the sensor library. Any ordered subset $S = \{\ell_1, \dots, \ell_K\} \subset \mathcal{L}$ of sensors can then form a (linear and continuous) *observation operator* through

$$L := [\ell_1, \dots, \ell_K]^T : \mathcal{U} \rightarrow \mathbb{R}^K, \quad Lu := [\ell_1(u), \dots, \ell_K(u)]^T.$$

The experimental measurements of L have the form

$$\mathbf{d} = [\ell_1(u) + \varepsilon_{\ell_1}, \dots, \ell_K(u) + \varepsilon_{\ell_K}]^T = Lu + \varepsilon \quad \text{with} \quad \varepsilon = [\varepsilon_{\ell_1}, \dots, \varepsilon_{\ell_K}]^T \sim \mathcal{N}(\mathbf{0}, \sigma^2 \Sigma_L), \quad (3)$$

where $\sigma^2 \Sigma_L$ is the noise covariance matrix defined through

$$\Sigma_L \in \mathbb{R}^{K \times K}, \quad \text{such that} \quad [\sigma^2 \Sigma_L]_{i,j} := \mathbf{cov}(\ell_j, \ell_i) = \mathbf{cov}(\varepsilon_{\ell_j}, \varepsilon_{\ell_i}) \quad (4)$$

with an auxiliary scaling parameter $\sigma^2 > 0$. We introduce σ^2 here as an additional variable to ease the discussion in the next section. However, we can set $\sigma^2 = 1$ without loss of generality (w.l.o.g.). We assume that the library \mathcal{L} and the noise covariance function \mathbf{cov} have been chosen such that Σ_L is s.p.d. for any combination of sensors in \mathcal{L} . This assumption gives rise to the L -dependent inner product and its induced norm

$$\langle \mathbf{d}, \tilde{\mathbf{d}} \rangle_{\Sigma_L^{-1}} := \mathbf{d}^T \Sigma_L^{-1} \tilde{\mathbf{d}}, \quad \|\mathbf{d}\|_{\Sigma_L^{-1}}^2 := \langle \mathbf{d}, \mathbf{d} \rangle_{\Sigma_L^{-1}}, \quad \forall \mathbf{d}, \tilde{\mathbf{d}} \in \mathbb{R}^K. \quad (5)$$

Measured with respect to this norm, the largest observation of any (normalized) state is thus

$$\gamma_L := \sup_{\|u\|_{\mathcal{U}}=1} \|Lu\|_{\Sigma_L^{-1}} = \sup_{u \in \mathcal{U}} \frac{\|Lu\|_{\Sigma_L^{-1}}}{\|u\|_{\mathcal{U}}}. \quad (6)$$

We show in Section 4.2 that γ_L increases under expansion of L with additional sensors despite the change in norm, and is therefore bounded by $\gamma_L \leq \gamma_{\mathcal{L}}$.

We also define the *parameter-to-observable map*

$$\mathbf{G}_{L,\theta} : \mathbb{R}^M \rightarrow \mathbb{R}^K, \quad \text{such that} \quad \mathbf{G}_{L,\theta} \mathbf{m} := Lu_\theta(\mathbf{m}). \quad (7)$$

With the assumptions above – in particular the linearity and uniform continuity (1) of u in \mathbf{m} – the map $\mathbf{G}_{L,\theta}$ is linear and uniformly bounded in \mathbf{m} . In a slight overload of notation, we identify the map $\mathbf{G}_{L,\theta} : \mathbb{R}^M \rightarrow \mathbb{R}^K$ with its (unique) matrix representation $\mathbf{G}_{L,\theta} \in \mathbb{R}^{K \times M}$ denote its matrix representation with respect to the unit basis $\{\mathbf{e}_m\}_{m=1}^M$. The likelihood of $\mathbf{d} \in \mathbb{R}^K$ obtained through the observation operator L for the parameter $\mathbf{m} \in \mathbb{R}^M$ and the system configuration θ is then

$$\Phi_L(\mathbf{d} \mid \mathbf{m}, \theta) := \frac{1}{\sqrt{2^K \det \Sigma_L}} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{d} - \mathbf{G}_{L,\theta} \mathbf{m}\|_{\Sigma_L^{-1}}^2\right).$$

Note that $\mathbf{G}_{L,\theta}$ may depend non-linearly on θ .

Posterior distribution

Once noisy measurement data $\mathbf{d} \approx Lu(\theta)$ is available, Bayes' theorem yields the posterior pdf as

$$\pi_{\text{post}}^{L,\theta}(\mathbf{m} \mid \mathbf{d}) = \frac{1}{Z(\theta)} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{G}_{L,\theta} \mathbf{m} - \mathbf{d}\|_{\Sigma_L^{-1}}^2 - \frac{1}{2} \|\mathbf{m} - \mathbf{m}_{\text{pr}}\|_{\Sigma_{\text{pr}}^{-1}}^2\right) \propto \pi_{\text{pr}}(\mathbf{m}) \cdot \Phi_L(\mathbf{d} \mid \mathbf{m}, \theta), \quad (8)$$

with normalization constant

$$Z(\theta) := \int_{\mathbb{R}^M} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{G}_{L,\theta} \mathbf{m} - \mathbf{d}\|_{\Sigma_L^{-1}}^2\right) d\mu_{\text{pr}}.$$

Due to the linearity of the parameter-to-observable map, the posterior measure $\mu_{\text{post}}^{L,\theta}$ is a Gaussian

$$\mu_{\text{post}}^{L,\theta} = \mathcal{N}(\mathbf{m}_{\text{post}}^{L,\theta}(\mathbf{d}), \Sigma_{\text{post}}^{L,\theta})$$

with known (cf. [1]) mean and covariance matrix

$$\mathbf{m}_{\text{post}}^{L,\theta}(\mathbf{d}) = \Sigma_{\text{post}}^{L,\theta} \left(\frac{1}{\sigma^2} \mathbf{G}_{L,\theta}^T \Sigma_L^{-1} \mathbf{d} + \Sigma_{\text{pr}}^{-1} \mathbf{m}_{\text{pr}} \right) \in \mathbb{R}^M, \quad (9)$$

$$\Sigma_{\text{post}}^{L,\theta} = \left(\frac{1}{\sigma^2} \mathbf{G}_{L,\theta}^T \Sigma_L^{-1} \mathbf{G}_{L,\theta} + \Sigma_{\text{pr}}^{-1} \right)^{-1} \in \mathbb{R}^{M \times M}. \quad (10)$$

The posterior $\mu_{\text{post}}^{L,\theta}$ thus depends not only on the choice of sensors, but also on the configuration θ under which their measurements were obtained. Therefore, to decrease the uncertainty in all possible posteriors with a single, θ -independent observation operator L , the construction of L should account for all admissible configurations $\theta \in \mathcal{P}$ under which u may be observed.

Remark 3. The linearity of $u_\theta(\mathbf{m})$ in \mathbf{m} is a strong assumption that dictates the Gaussian posterior. However, in combination with the hyper-parameter θ , our setting here can be re-interpreted as the Laplace-approximation (cf., [51]) for a non-linear state map $\theta \mapsto u(\theta)$ with uncertain $\theta \in \mathbb{R}^p$. In this case, the hyper-parameter set $\mathcal{P} \subset \mathbb{R}^p$ here encloses those θ that, according to their prior, are potential maximum a posteriori probability (MAP) points of the non-linear inverse problem, while the linear parameter \mathbf{m} becomes a scaling vector in the linearization of $\theta \mapsto u(\theta)$ around $\theta \in \mathcal{P}$. A more detailed interpretation of our setting here in terms of the Laplace approximation of non-linear inverse problems is provided in [52]. However, the additional analysis and numerical treatment of this setting is beyond the scope of this work and part of ongoing research.

3. The observability coefficient

In this section, we characterize how the choice of sensors in the observation operator L and its associated noise covariance matrix Σ_L influence the uncertainty in the posteriors $\mu_{\text{post}}^{L,\theta}$, $\theta \in \mathcal{P}$. We identify an observability coefficient that bounds the eigenvalues of the posterior covariance matrices $\Sigma_{\text{post}}^{L,\theta}$, $\theta \in \mathcal{P}$ with respect to L , and facilitates the sensor selection algorithm presented in Section 4.

3.1. Eigenvalues of the posterior covariance matrix

The uncertainty in the posterior $\pi_{\text{post}}^{L,\theta}$ for any configuration $\theta \in \mathcal{P}$ is uniquely characterized by the posterior covariance matrix $\Sigma_{\text{post}}^{L,\theta}$, which is in turn connected to the observation operator L through the parameter-to-observable map $\mathbf{G}_{L,\theta}$ and the noise covariance matrix Σ_L . To measure the uncertainty in $\Sigma_{\text{post}}^{L,\theta}$, the OED literature suggests a variety of different utility functions to be minimized over L in order to optimize the sensor choice. Many of these utility functions can be expressed in terms of the eigenvalues $\lambda_L^{\theta,1} \geq \dots \geq \lambda_L^{\theta,M} > 0$ of $\Sigma_{\text{post}}^{L,\theta}$, e.g.,

$$\begin{aligned} \text{A-OED:} \quad & \text{trace}(\Sigma_{\text{post}}^{L,\theta}) = \sum_{m=1}^M \lambda_L^{\theta,m} & (\text{mean variance}) \\ \text{D-OED:} \quad & \det(\Sigma_{\text{post}}^{L,\theta}) = \prod_{m=1}^M \lambda_L^{\theta,m} & (\text{volume}) \\ \text{E-OED:} \quad & \lambda_{\max}(\Sigma_{\text{post}}^{L,\theta}) = \lambda_L^{\theta,1} & (\text{spectral radius}). \end{aligned}$$

In practice, the choice of the utility function is dictated by the application. In E-OED, for instance, posteriors whose uncertainty ellipsoids stretch out into any one direction are avoided, whereas D-OED minimizes the overall volume of the uncertainty ellipsoid regardless of the uncertainty in any one parameter direction. We refer to [5] for a detailed introduction and other OED criteria.

Considering the hyper-parameterized setting where each configuration θ influences the posterior uncertainty, we seek to choose a single observation operator L such that the selected utility function remains small for *all* configurations $\theta \in \mathcal{P}$, e.g., for E-OED, minimizing

$$\min_{\ell_1, \dots, \ell_K \in \mathcal{L}} \max_{\theta \in \mathcal{P}} \lambda_{\max}(\Sigma_{\text{post}}^{L,\theta}) \quad \text{such that} \quad L = [\ell_1, \dots, \ell_K]^T$$

guarantees that the longest axis of each posterior covariance matrix $\Sigma_{\text{post}}^{L,\theta}$ for any $\theta \in \mathcal{P}$ has the same guaranteed upper bound. The difficulty here is that the minimization over \mathcal{P} necessitates repeated, cost-intensive model evaluations to compute the utility function for many different configurations θ . In the following, we therefore introduce an upper bound to the posterior eigenvalues that can be optimized through an observability criterion with far fewer model solves. The bound's optimization indirectly reduces the different utility functions through the posterior eigenvalues.

Recalling that $\Sigma_{\text{post}}^{L,\theta}$ is s.p.d., let $\{\psi_m\}_{m=1}^M$ be an orthonormal eigenvector basis of $\Sigma_{\text{post}}^{L,\theta}$, i.e. $\psi_m^T \psi_n = \delta_{m,n}$ and

$$\Sigma_{\text{post}}^{L,\theta} \psi_m = \lambda_L^{\theta,m} \psi_m \quad m = 1, \dots, M. \quad (11)$$

Using the representation (10), any eigenvalue $\lambda_L^{\theta,m}$ can be written in the form

$$\frac{1}{\lambda_L^{\theta,m}} = \psi_m^T \left[\Sigma_{\text{post}}^{L,\theta} \right]^{-1} \psi_m = \psi_m^T \left[\frac{1}{\sigma^2} \mathbf{G}_{L,\theta}^T \Sigma_L^{-1} \mathbf{G}_{L,\theta} + \Sigma_{\text{pr}}^{-1} \right] \psi_m = \frac{1}{\sigma^2} \|\mathbf{G}_{L,\theta} \psi_m\|_{\Sigma_L^{-1}}^2 + \|\psi_m\|_{\Sigma_{\text{pr}}^{-1}}^2. \quad (12)$$

Since ψ_m depends implicitly on L and θ through (11), we cannot use this representation directly to optimize over L . To take out the dependency on ψ_m , we bound $\|\psi_m\|_{\Sigma_{\text{pr}}^{-1}}^2 \geq \frac{1}{\lambda_{\text{pr}}^{\max}}$ in terms of the maximum eigenvalue of the prior covariance matrix Σ_{pr} . Likewise, we define the Rayleigh quotient

$$\beta_G(\theta, L) := \inf_{\mathbf{m} \in \mathbb{R}^M} \frac{\|\mathbf{G}_{L,\theta} \mathbf{m}\|_{\Sigma_L^{-1}}}{\|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}}} = \inf_{\mathbf{m} \in \mathbb{R}^M} \frac{\|Lu_\theta(\mathbf{m})\|_{\Sigma_L^{-1}}}{\|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}}}, \quad (13)$$

as the minimum ratio between an observation for a parameter \mathbf{m} relative to the prior's covariance norm. We call $\beta_G(\theta, L)$ *observability coefficient* in reference to optimal control theory (cf., [53], chapter 1.2). The maximization of $\beta_G(\theta, L)$ has been explored for uncorrelated noise in the E-OED literature in [16]. In this work, θ was treated as a nuisance parameter, with the optimization problem for L posed over a submatrix of the Fisher information matrix. When disregarding the hyper-parameter θ , $\mathbf{G}_{L,\theta}$ is comparable to the *inf-sup* stability constant in the PBDW method [54,20], though a θ -dependent variation was introduced in [18] in the context of 3D-VAR variational data assimilation. The maximization of the PBDW *inf-sup* stability constant is typically performed in a greedy OMP procedure, with convergence properties analyzed in [19,55]. However, while the PBDW noise model is indeed correlated when translated in the Bayesian setting, the noise covariance matrix Σ_L has a very specific structure defined through the Riesz representation of the chosen sensors (see [52], section 3.4.1 for the connection). In this sense, $\beta_G(\theta, L)$ can be considered a generalization of the PBDW *inf-sup* stability constant to arbitrary correlated noise models.

From (12) and (13) we obtain the upper bound

$$\lambda_L^{\theta,m} = \left(\frac{1}{\sigma^2} \frac{\|\mathbf{G}_{L,\theta} \psi_m\|_{\Sigma_L^{-1}}^2}{\|\psi_m\|_{\Sigma_{\text{pr}}^{-1}}^2} + 1 \right)^{-1} \|\psi_m\|_{\Sigma_{\text{pr}}^{-1}}^{-2} \leq \left(\frac{1}{\sigma^2} \beta_G(\theta, L)^2 + 1 \right)^{-1} \lambda_{\text{pr}}^{\max}. \quad (14)$$

Geometrically, this bound means that the radius $\lambda_L^{\theta,1}$ of the outer ball around the posterior uncertainty ellipsoid is smaller than that of the prior uncertainty ellipsoid by at least the factor $\left(\frac{1}{\sigma^2} \beta_G(\theta, L)^2 + 1 \right)^{-1}$. As expected, the influence of L is strongest when the measurement noise is small such that data can be trusted ($\sigma^2 \ll 1$), and diminishes with increasing noise levels ($\sigma^2 \gg 1$).

The main idea of our sensor selection method outlined in Section 4 is to choose L to maximize $\min_\theta \beta_G(\theta, L)$, which, geometrically, corresponds to minimizing the outer ball with radius $\max_\theta \lambda_L^{\theta,1}$ containing all uncertainty ellipsoids (i.e., for any $\theta \in \mathcal{P}$). By definition of the A-OED, D-OED, and E-OED utility functions, this approach is most fitting for E-OED where only the largest eigenvalue is measured in the utility function. In particular, if the prior is independent and identically distributed, i.e. Σ_{pr} is a multiplication of the identity matrix, then the upper bound (14) is indeed equal to the maximum posterior eigenvalue $\lambda_L^{\theta,1}$. For A-OED the maximization of $\min_\theta \beta_G(\theta, L)$ remains applicable when the parameter dimension is low or the eigenvalues of the posterior covariance matrices decay slowly such that the A-OED criterion is dominated by the most uncertain parameter directions. In contrast, the approach is less suitable for D-OED which is more sensitive to the relative improvement in each parameter direction rather than the worst case direction. We illustrate the correlation between these three OED criteria and the observability coefficient on a practical example in Section 5, in particular Figs. 4, 5, 6.

3.2. Parameter restriction

An essential property of $\beta_G(\theta, L)$ is that $\beta_G(\theta, L) = 0$ if $K < M$, i.e., the number of sensors in L is smaller than the number of parameter dimensions. In this case, $\beta_G(\theta, L)$ cannot distinguish between sensors during the first $M - 1$ steps of an iterative algorithm, or in general when less than a total of M sensors are supposed to be chosen. For medium-dimensional parameter spaces (M in the order of tens), we mitigate this issue by restricting \mathbf{m} to the subspace $\text{span}\{\varphi_1, \dots, \varphi_{\min\{K, M\}}\} \subset \mathbb{R}^M$ spanned by the first $\min\{K, M\}$ eigenvectors of Σ_{pr} corresponding to its largest eigenvalues, i.e., the subspace with the largest prior uncertainty. For high-dimensional parameter spaces or when the model \mathcal{M}_θ has a non-trivial null-space, we bound $\beta_G(\theta, L)$ further

$$\beta_G(\theta, L) = \inf_{\mathbf{m} \in \mathbb{R}^M} \frac{\|Lu_\theta(\mathbf{m})\|_{\Sigma_L^{-1}}}{\|u_\theta(\mathbf{m})\|_{\mathcal{V}}} \frac{\|u_\theta(\mathbf{m})\|_{\mathcal{V}}}{\|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}}} \geq \inf_{x \in \mathcal{V}_\theta} \frac{\|Lx\|_{\Sigma_L^{-1}}}{\|x\|_{\mathcal{V}}} \inf_{\mathbf{m} \in \mathbb{R}^M} \frac{\|u_\theta(\mathbf{m})\|_{\mathcal{V}}}{\|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}}} = \beta_V(\theta, L) \underline{\eta}(\theta) \quad (15)$$

where we define the linear space \mathcal{V}_θ of all achievable states

$$\mathcal{V}_\theta := \{u_\theta(\mathbf{m}) \in \mathcal{V} : \mathbf{m} \in \mathbb{R}^M\}$$

and the coefficients

$$\beta_V(\theta, L) := \inf_{u \in \mathcal{V}_\theta} \frac{\|Lu\|_{\Sigma_L^{-1}}}{\|u\|_{\mathcal{V}}}, \quad \underline{\eta}(\theta) := \inf_{\mathbf{m} \in \mathbb{R}^M} \frac{\|u_\theta(\mathbf{m})\|_{\mathcal{V}}}{\|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}}}. \quad (16)$$

The value of $\underline{\eta}(\theta)$ describes the minimal state change that a parameter \mathbf{m} can achieve relative to its prior-induced norm $\|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}}$. It can filter out parameter directions that have little influence on the states $u_\theta(\mathbf{m})$. In contrast, the observability coefficient $\beta_V(\theta, L)$ depends on the prior only implicitly via \mathcal{V}_θ ; it quantifies the minimum amount of information (measured with respect to the noise model) that can be obtained on any state in \mathcal{V}_θ relative to its norm. Future work will investigate how to optimally restrict the parameter space based on $\underline{\eta}(\theta)$ before choosing sensors that maximize $\beta_V(\theta, L)$. Existing parameter reduction approaches in a similar context include [37,48,49,36]. In this work, however, we solely focus on the maximization of $\beta_G(\theta, L)$ and, by extension, $\beta_V(\theta, L)$ and henceforth assume that M is sufficiently small and $\underline{\eta} := \inf_{\theta \in \mathcal{P}} \underline{\eta}(\theta) > 0$ is bounded away from zero.

3.3. Observability under model approximations

To optimize the observability coefficient $\beta_G(\theta, L)$ or $\beta_V(\theta, L)$, it must be computed for many different configurations $\theta \in \mathcal{P}$. The accumulating computational cost motivates the use of *reduced-order* surrogate models, which typically yield considerable computational savings versus the original *full-order* model. However, this leads to errors in the state approximation. In the following, we thus quantify the influence of state approximation error on the observability coefficients $\beta_G(\theta, L)$ and $\beta_V(\theta, L)$. An analysis of the change in posterior distributions when the entire model \mathcal{M}_θ is substituted in the inverse problem can be found in [1], section 4.4 (pp. 504–508), and the references therein.

Suppose a reduced-order surrogate model $\tilde{\mathcal{M}}_\theta(\tilde{u}_\theta(\mathbf{m}); \mathbf{m}) = 0$ is available that yields for any configuration $\theta \in \mathcal{P}$ and parameter $\mathbf{m} \in \mathbb{R}^M$ a unique solution $\tilde{u}_\theta(\mathbf{m}) \in \mathcal{U}$ such that

$$\|u_\theta(\mathbf{m}) - \tilde{u}_\theta(\mathbf{m})\|_{\mathcal{U}} \leq \varepsilon_\theta \|u_\theta(\mathbf{m})\|_{\mathcal{U}} \quad \text{with accuracy} \quad 0 \leq \varepsilon_\theta \leq \varepsilon < 1. \quad (17)$$

Analogously to (13) and (16), we define the reduced-order observability coefficients

$$\tilde{\beta}_G(\theta, L) := \inf_{\mathbf{m} \in \mathbb{R}^M} \frac{\|L\tilde{u}_\theta(\mathbf{m})\|_{\Sigma_L^{-1}}}{\|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}}}, \quad \tilde{\beta}_V(\theta, L) := \inf_{\mathbf{m} \in \mathbb{R}^M} \frac{\|L\tilde{u}_\theta(\mathbf{m})\|_{\Sigma_L^{-1}}}{\|\tilde{u}_\theta(\mathbf{m})\|_{\mathcal{U}}} \quad (18)$$

to quantify the smallest observations of the surrogate states. For many applications, it is possible to choose a reduced-order model whose solution can be computed at a significantly reduced cost such that $\tilde{\beta}_G(\theta, L)$ and $\tilde{\beta}_V(\theta, L)$ are much cheaper to compute than their full-order counterparts $\beta_G(\theta, L)$ and $\beta_V(\theta, L)$. Since the construction of such a surrogate model depends strongly on the application itself, we refer to the literature (e.g., [56–60]) for tangible approaches.

Recalling the definition of γ_L in (6), we start by bounding how closely the surrogate observability coefficient $\tilde{\beta}_V(\theta, L)$ approximates the full-order $\beta_V(\theta, L)$.

Proposition 1. Let $\eta(\theta) > 0$ hold, and let $\tilde{u}_\theta(\mathbf{m}) \in \mathcal{U}$ be an approximation to $u_\theta(\mathbf{m})$ such that (17) holds for all $\theta \in \mathcal{P}$, $\mathbf{m} \in \mathbb{R}^M$. Then

$$(1 - \varepsilon_\theta) \tilde{\beta}_V(\theta, L) - \gamma_L \varepsilon_\theta \leq \beta_V(\theta, L) \leq (1 + \varepsilon_\theta) \tilde{\beta}_V(\theta, L) + \gamma_L \varepsilon_\theta. \quad (19)$$

Proof. Let $\mathbf{m} \in \mathbb{R}^M \setminus \{\mathbf{0}\}$ be arbitrary. Using (17) and the (reversed) triangle inequality, we obtain the bound

$$\frac{\|\tilde{u}_\theta(\mathbf{m})\|_{\mathcal{U}}}{\|u_\theta(\mathbf{m})\|_{\mathcal{U}}} \geq \frac{\|u_\theta(\mathbf{m})\|_{\mathcal{U}} - \|u_\theta(\mathbf{m}) - \tilde{u}_\theta(\mathbf{m})\|_{\mathcal{U}}}{\|u_\theta(\mathbf{m})\|_{\mathcal{U}}} \geq 1 - \varepsilon_\theta. \quad (20)$$

Note here that $\eta(\theta) > 0$ implies $\|u_\theta(\mathbf{m})\|_{\mathcal{U}} > 0$ so the quotient is indeed well defined. The ratio of observation to state can now be bounded from below by

$$\begin{aligned} \frac{\|Lu_\theta(\mathbf{m})\|_{\Sigma_L^{-1}}}{\|u_\theta(\mathbf{m})\|_{\mathcal{U}}} &\geq \frac{\|L\tilde{u}_\theta(\mathbf{m})\|_{\Sigma_L^{-1}}}{\|u_\theta(\mathbf{m})\|_{\mathcal{U}}} - \frac{\|L(u_\theta(\mathbf{m}) - \tilde{u}_\theta(\mathbf{m}))\|_{\Sigma_L^{-1}}}{\|u_\theta(\mathbf{m})\|_{\mathcal{U}}} \\ &\geq \frac{\|\tilde{u}_\theta(\mathbf{m})\|_{\mathcal{U}}}{\|u_\theta(\mathbf{m})\|_{\mathcal{U}}} \frac{\|L\tilde{u}_\theta(\mathbf{m})\|_{\Sigma_L^{-1}}}{\|\tilde{u}_\theta(\mathbf{m})\|_{\mathcal{U}}} - \gamma_L \frac{\|u_\theta(\mathbf{m}) - \tilde{u}_\theta(\mathbf{m})\|_{\mathcal{U}}}{\|u_\theta(\mathbf{m})\|_{\mathcal{U}}} \\ &\geq (1 - \varepsilon_\theta) \frac{\|L\tilde{u}_\theta(\mathbf{m})\|_{\Sigma_L^{-1}}}{\|\tilde{u}_\theta(\mathbf{m})\|_{\mathcal{U}}} - \gamma_L \varepsilon_\theta \\ &\geq (1 - \varepsilon_\theta) \tilde{\beta}_V(\theta, L) - \gamma_L \varepsilon_\theta, \end{aligned}$$

where we have applied the reverse triangle inequality, definition (6), the bounds (17), (20), and definition (18) of $\tilde{\beta}_V(\theta, L)$. Since \mathbf{m} is arbitrary, the lower bound in (19) follows from definition (13) of $\beta_V(\theta, L)$. The upper bound in (19) follows analogously. \square

For the observability of the parameter-to-observable map $\mathbf{G}_{L,\theta}$ and its approximation $\mathbf{m} \mapsto L\tilde{u}_\theta(\mathbf{m})$, we obtain a similar bound. It uses the norm $\tilde{\eta}(\theta)$ of $u_\theta : \mathbf{m} \mapsto u_\theta(\mathbf{m})$ as a map from the parameter to the state space, see (1).

Proposition 2. Let $\tilde{u}_\theta(\mathbf{m}) \in \mathcal{U}$ be an approximation to $u_\theta(\mathbf{m})$ such that (17) holds for all $\theta \in \mathcal{P}$, $\mathbf{m} \in \mathbb{R}^M$. Then

$$\tilde{\beta}_G(\theta, L) - \gamma_L \tilde{\eta}(\theta) \varepsilon_\theta \leq \beta_G(\theta, L) \leq \tilde{\beta}_G(\theta, L) + \gamma_L \tilde{\eta}(\theta) \varepsilon_\theta. \quad (21)$$

Proof. Let $\mathbf{m} \in \mathbb{R}^M \setminus \{\mathbf{0}\}$ be arbitrary. Then

$$\begin{aligned} \|Lu_\theta(\mathbf{m})\|_{\Sigma_L^{-1}} &\geq \|L\tilde{u}_\theta(\mathbf{m})\|_{\Sigma_L^{-1}} - \|L(u_\theta(\mathbf{m}) - \tilde{u}_\theta(\mathbf{m}))\|_{\Sigma_L^{-1}} \\ &\geq \|L\tilde{u}_\theta(\mathbf{m})\|_{\Sigma_L^{-1}} - \gamma_L \|u_\theta(\mathbf{m}) - \tilde{u}_\theta(\mathbf{m})\|_{\mathcal{U}} \\ &\geq \|L\tilde{u}_\theta(\mathbf{m})\|_{\Sigma_L^{-1}} - \gamma_L \varepsilon_\theta \|u_\theta(\mathbf{m})\|_{\mathcal{U}} \end{aligned}$$

$$\geq \|L\tilde{u}_\theta(\mathbf{m})\|_{\Sigma_L^{-1}} - \gamma_L \epsilon_\theta \bar{\eta}(\theta) \|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}},$$

where we have used the reverse triangle inequality, followed by (6), (17), and (1). We divide by $\|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}}$ and take the infimum over \mathbf{m} to obtain

$$\beta_G(\theta, L) = \inf_{\mathbf{m} \in \mathbb{R}^M} \frac{\|Lu_\theta(\mathbf{m})\|_{\Sigma_L^{-1}}}{\|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}}} \geq \inf_{\mathbf{m} \in \mathbb{R}^M} \frac{\|L\tilde{u}_\theta(\mathbf{m})\|_{\Sigma_L^{-1}}}{\|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}}} - \gamma_L \bar{\eta}(\theta) \epsilon_\theta = \tilde{\beta}_G(\theta, L) - \gamma_L \bar{\eta}(\theta) \epsilon_\theta.$$

The upper bound in (21) follows analogously. \square

If ϵ_θ is sufficiently small, Propositions 1 and 2 justify employing the surrogates $\tilde{\beta}_V(\theta, L)$ and $\tilde{\beta}_G(\theta, L)$ instead of the original full-order observability coefficients $\beta_V(\theta, L)$ and $\beta_G(\theta, L)$. This substitution becomes especially necessary when the computation of $u_\theta(\mathbf{m})$ is too expensive to evaluate $\beta_V(\theta, L)$ or $\beta_G(\theta, L)$ repeatedly for different configurations θ .

Another approximation step in our sensor selection algorithm relies on the identification of a parameter direction $\mathbf{v} \in \mathbb{R}^M$ with comparatively small observability, i.e.

$$\frac{\|Lu_\theta(\mathbf{v})\|_{\Sigma_L^{-1}}}{\|\mathbf{v}\|_{\Sigma_{\text{pr}}^{-1}}} \approx \inf_{\mathbf{m} \in \mathbb{R}^M} \frac{\|Lu_\theta(\mathbf{m})\|_{\Sigma_L^{-1}}}{\|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}}} = \beta_G(\theta, L) \quad \text{or} \quad \frac{\|Lu_\theta(\mathbf{v})\|_{\Sigma_L^{-1}}}{\|u_\theta(\mathbf{v})\|_{\mathcal{U}}} \approx \inf_{u \in \mathcal{U}} \frac{\|Lu\|_{\Sigma_L^{-1}}}{\|u\|_{\mathcal{U}}} = \beta_V(\theta, L).$$

The ideal choice would be the infimizer of respectively $\beta_G(\theta, L)$ or $\beta_V(\theta, L)$, but its computation involves M full-order model evaluations (cf. Section 4.4). To avoid these costly computations, we instead choose \mathbf{v} as the infimizer of the respective reduced-order observability coefficient. This choice is justified for small $\epsilon_\theta < 1$ by the following proposition:

Proposition 3. Let $\bar{\eta}(\theta) > 0$ hold, and let $\tilde{u}_\theta(\mathbf{m}) \in \mathcal{U}$ be an approximation to $u_\theta(\mathbf{m})$ such that (17) holds for all $\theta \in \mathcal{P}$, $\mathbf{m} \in \mathbb{R}^M$. Suppose $\mathbf{v} \in \arg \inf_{\mathbf{m} \in \mathbb{R}^M} \|L\tilde{u}_\theta(\mathbf{m})\|_{\Sigma_L^{-1}} \|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}}^{-1}$, then

$$\beta_G(\theta, L) \leq \frac{\|Lu_\theta(\mathbf{v})\|_{\Sigma_L^{-1}}}{\|\mathbf{v}\|_{\Sigma_{\text{pr}}^{-1}}} \leq \beta_G(\theta, L) + 2\gamma_L \bar{\eta}(\theta) \epsilon_\theta. \quad (22)$$

Likewise, if $\mathbf{v} \in \arg \inf_{\mathbf{m} \in \mathbb{R}^M} \|\tilde{u}_\theta(\mathbf{m})\|_{\mathcal{U}}^{-1} \|L\tilde{u}_\theta(\mathbf{m})\|_{\Sigma_L^{-1}}$, then

$$\beta_V(\theta, L) \leq \frac{\|Lu_\theta(\mathbf{v})\|_{\Sigma_L^{-1}}}{\|u_\theta(\mathbf{v})\|_{\mathcal{U}}} \leq \frac{1 + \epsilon_\theta}{1 - \epsilon_\theta} (\beta_V(\theta, L) + \gamma_L \epsilon_\theta) + \gamma_L \epsilon_\theta. \quad (23)$$

Proof. For both (22) and (23) the lower bound follows directly from definitions (13) and (16). To prove the upper bound in (22), let $\mathbf{v} \in \arg \inf_{\mathbf{m} \in \mathbb{R}^M} \|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}}^{-1} \|L\tilde{u}_\theta(\mathbf{m})\|_{\Sigma_L^{-1}}$. Following the same steps as in the proof of Proposition 2, we can then bound

$$\frac{\|Lu_\theta(\mathbf{v})\|_{\Sigma_L^{-1}}}{\|\mathbf{v}\|_{\Sigma_{\text{pr}}^{-1}}} \leq \frac{\|L\tilde{u}_\theta(\mathbf{v})\|_{\Sigma_L^{-1}}}{\|\mathbf{v}\|_{\Sigma_{\text{pr}}^{-1}}} + \frac{\|L(u_\theta(\mathbf{v}) - \tilde{u}_\theta(\mathbf{v}))\|_{\Sigma_L^{-1}}}{\|\mathbf{v}\|_{\Sigma_{\text{pr}}^{-1}}} \leq \tilde{\beta}_G(\theta, L) + \gamma_L \bar{\eta}(\theta) \epsilon_\theta.$$

The upper bound in (22) then follows with Proposition 2.

To prove the upper bound in (23), let $\mathbf{v} \in \arg \inf_{\mathbf{m} \in \mathbb{R}^M} \|\tilde{u}_\theta(\mathbf{m})\|_{\mathcal{U}}^{-1} \|L\tilde{u}_\theta(\mathbf{m})\|_{\Sigma_L^{-1}}$. Then

$$\frac{\|Lu_\theta(\mathbf{v})\|_{\Sigma_L^{-1}}}{\|u_\theta(\mathbf{v})\|_{\mathcal{U}}} \leq \frac{\|L\tilde{u}_\theta(\mathbf{v})\|_{\Sigma_L^{-1}}}{\|\tilde{u}_\theta(\mathbf{v})\|_{\mathcal{U}}} \frac{\|\tilde{u}_\theta(\mathbf{v})\|_{\mathcal{U}}}{\|u_\theta(\mathbf{v})\|_{\mathcal{U}}} + \frac{\|L(u_\theta(\mathbf{v}) - \tilde{u}_\theta(\mathbf{v}))\|_{\Sigma_L^{-1}}}{\|u_\theta(\mathbf{v})\|_{\mathcal{U}}} \leq (1 + \epsilon) \tilde{\beta}_V(\theta, L) + \gamma_L \epsilon_\theta.$$

The result then follows with Proposition 1. \square

4. Sensor selection

In this section, we present a sensor selection algorithm that iteratively chooses the individual sensors in L to increase the minimal observability coefficient $\min_{\theta \in \mathcal{P}} \beta_G(\theta, L)$ and thereby decreases the upper bound (14) for the eigenvalues of the posterior covariance matrix for all admissible system configurations $\theta \in \mathcal{P}$. The advantage of using the observability coefficient rather than targeting a utility function directly is that neither the posterior covariance matrix nor its action need to be computed to evaluate the improvement that any additional sensor would bring. Compared to computing the posterior, using the observability coefficient $\beta_G(\theta, L)$ as target for the iterative “max-min” optimization reduces the number of full-order model solves per iteration from M to a single one. Although the iterative procedure cannot guarantee finding the optimal observability over all sensor combinations, the underlying greedy searches are well-established in practice, and can be shown to perform with exponentially decreasing error rates in closely related settings, see [61,19,55,62,63]. The iterative approach is also relatively easy to implement, allows a simple way of dealing with combinatorial restrictions, and can deal with large sensor libraries even for correlated noise.

Algorithm 1: SensorSelection.

Input: sensor library $\mathcal{L} \subset \mathcal{U}'$, training set $\Xi_{\text{train}} \subset \mathcal{P}$, maximum number of sensors $K_{\text{max}} \leq |\mathcal{L}|$, $K_{\text{max}} \geq M$, surrogate model $\tilde{\mathcal{M}}_\theta$, covariance function $\text{cov} : \mathcal{L} \times \mathcal{L} \rightarrow \mathbb{R}$

Compute $\Sigma_{\text{pr}} = [\varphi_1, \dots, \varphi_M]^T \mathbf{D}_{\text{pr}} [\varphi_1, \dots, \varphi_M]$ // eigenvalue decomposition

For all $\theta \in \Xi_{\text{train}}$, $1 \leq m \leq M$, compute $\tilde{u}_\theta(\varphi_m)$ // preparation

$K \leftarrow 0$, $\theta_0 \leftarrow \arg \max_{\theta \in \Xi_{\text{train}}} \|\tilde{u}_\theta(\varphi_1)\|_{L^1}$, $\mathbf{m}_0 \leftarrow \varphi_1$ // initialization

while $K < K_{\text{max}}$ **do**

// data matching

Solve full-order equation $\mathcal{M}_{\theta_K}(u_K, \mathbf{m}_K)$ for u_K // "worst-case" state

$\ell_{K+1} \leftarrow \arg \max_{\ell \in \mathcal{L}} \text{ObservabilityGain}(L, \mathbf{C}_L, \ell)$ // sensor selection

$L, \Sigma_L, \mathbf{C}_L \leftarrow \text{CholeskyExpansion}(L, \Sigma_L, \mathbf{C}_L, \ell_{K+1})$ // expansion

$K \leftarrow K + 1$

// greedy search

for $\theta \in \Xi_{\text{train}}$ **do**

$\tilde{\beta}_V(\theta, L), \mathbf{m}_{\min}(\theta) \leftarrow \text{SurrogateObservability}(\theta, L, \mathbf{C}_L)$ // update coefficients

$\theta_K \leftarrow \arg \min_{\theta \in \Xi_{\text{train}}} \tilde{\beta}_V(\theta, L)$ // "worst-case" hyper-parameter

$\mathbf{m}_K \leftarrow \sum_{m=1}^{\min(M, K)} [\mathbf{m}_{\min}(\theta_K)]_m \varphi_m$ // "worst-case" parameter

return L, \mathbf{C}_L

We start with a high-level overview of the algorithm to describe its main steps and ideas. To keep the exposition simple, we delay introducing computational details on three key operations – the matrix expansions, the computations of the observability gain, and the surrogate observability coefficient – to Sections 4.2, 4.3 and 4.4. In the algorithm description, these operations are denoted as function calls `CholeskyExpansion`, `ObservabilityGain`, `SurrogateObservability`, to be presented in Algorithms 2, 3, and 4 later.

4.1. Sensor selection

Our goal is to identify an observation operator L composed of K_{max} sensors from the sensor library \mathcal{L} that solves the “max-min”-optimization problem

$$\max_L \min_{\theta \in \mathcal{P}} \beta_G(\theta, L) \quad \text{such that } L = [\ell_1, \dots, \ell_{K_{\text{max}}}]^T : \mathcal{U} \rightarrow \mathbb{R}^{K_{\text{max}}} \text{ with } \ell_1, \dots, \ell_{K_{\text{max}}} \in \mathcal{L}.$$

In our sensor selection algorithm, we iteratively expand the observation operator L . It will be shown in Section 4.3 that $\beta_G(\theta, L)$ is monotonously increasing under expansion of L , which guarantees that this iterative expansion will indeed increase the observability coefficient $\beta_G(\theta, L)$ for all $\theta \in \mathcal{P}$. In each iteration, the algorithm performs two main steps:

- A **greedy search** over a training set $\Xi_{\text{train}} \subset \mathcal{P}$ to identify the configuration $\theta \in \Xi_{\text{train}}$ for which the observability coefficient $\beta_G(\theta, L)$ is (approximately) minimal. At this “worst-case” configuration, we also identify the corresponding “worst-case” parameter $\mathbf{m} = \arg \min_{\mathbf{m} \in \mathbb{R}^M} \|Lu_\theta(\mathbf{m})\|_{\Sigma_L^{-1}} \|\mathbf{m}\|_{\Sigma_{\text{pr}}^{-1}}^{-1}$.
- A **data-matching step** to identify a sensor in the library \mathcal{L} that maximizes the observation of the “worst-case” parameter at the “worst-case” configuration θ .

To keep the computational effort feasible, we use a reduced-order model to approximate $\beta_G(\theta, L)$ during the greedy search. The procedure is summarized in Algorithm 1, using forward references `CholeskyExpansion` (Algorithm 2, Section 4.2), `ObservabilityGain` (Algorithm 3, Section 4.3), `SurrogateObservability` (Algorithm 4, Section 4.4) to the explanations of the more involved sub-steps. In the following, we explain the high-level computational details.

Preparations

Let $\Sigma_{\text{pr}} = \mathbf{U}^T \mathbf{D}_{\text{pr}} \mathbf{U}$ be the eigenvalue decomposition of the s.p.d. prior covariance matrix with $\mathbf{U} = [\varphi_1, \dots, \varphi_M] \in \mathbb{R}^{M \times M}$, $\varphi_j \in \mathbb{R}^M$ orthonormal in the Euclidean inner product, and $\mathbf{D}_{\text{pr}} = \text{diag}(\lambda_{\text{pr}}^1, \dots, \lambda_{\text{pr}}^M)$ a diagonal matrix containing the eigenvalues $\lambda_{\text{pr}}^1 \geq \dots \geq \lambda_{\text{pr}}^M > 0$ in decreasing order. As discussed in Section 3.2, we assume M is sufficiently small for $\beta_G(\theta, L)$ to be meaningful when all K_{max} sensors have been chosen, i.e., $M \leq K_{\text{max}}$. In order to increase $\beta_G(\theta, L)$ uniformly throughout the hyper-parameter domain \mathcal{P} , we choose a finite training set, $\Xi_{\text{train}} \subset \mathcal{P}$, that is fine enough to capture the θ -dependent variations in the state $u_\theta(\mathbf{m})$. We assume a reduced-order model is available such that we can compute approximations $\tilde{u}_\theta(\varphi_m) \approx u_\theta(\varphi_m)$ for each $\theta \in \Xi_{\text{train}}$ and $1 \leq m \leq M$ within an acceptable computation time while also guaranteeing the accuracy requirement (17). If necessary, the two criteria can be balanced via adaptive training domains (e.g., [64,65]). The reduced-order model will be used in each iteration to identify the hyper-parameter θ with the worst observability for the current observation operator L . If memory allows (e.g., with projection-based

surrogate models), the surrogate states $\tilde{u}_\theta(\varphi_m)$ for $\theta \in \Xi_{\text{train}}$, $1 \leq m \leq M$ should be computed and stored at the start of the algorithm to avoid re-computations.

As a first “worst-case” parameter direction, \mathbf{m}_0 , we choose the vector φ_1 with the largest prior uncertainty. Likewise, we choose the “worst-case” configuration $\theta_K \in \Xi_{\text{train}}$ as the one for which the corresponding state $\tilde{u}_\theta(\varphi_1)$ is the largest.

Data-matching step

In each iteration, we first compute the full-order state $u_K = u_{\theta_K}(\mathbf{m}_K)$ at the “worst-case” parameter \mathbf{m}_K and configuration θ_K . We use the full-order state $u_{\theta_K}(\mathbf{m}_K)$ rather than its reduced-order surrogate in order to avoid training on local approximation inaccuracies in the reduced-order model. However, since we only require a single full-order model solve per iteration, we are still keeping the computational effort small compared to the M full-order model solves that would be required for setting up the entire posterior covariance matrix $\Sigma_{\text{post}}^{L,\theta}$.

We continue by computing for each sensor ℓ in the sensor library \mathcal{L} by how much the norm $\|L\mathbf{m}_K\|_{\Sigma_L^{-1}}$ would increase if L was expanded to $[L^T, \ell]^T$, i.e. we compute the *observability gain*

$$\| [L, \ell](u) \|_{\Sigma_{[L, \ell]}^{-1}}^2 - \| Lu \|_{\Sigma_L^{-1}}^2.$$

This value is returned by the function `ObservabilityGain` defined in Algorithm 3 below. It will be shown in Section 4.3 that it can be implemented efficiently without explicitly computing $\| [L, \ell](u) \|_{\Sigma_{[L, \ell]}^{-1}}^2$ or $\| Lu \|_{\Sigma_L^{-1}}^2$.

We choose the sensor ℓ_{K+1} as the one with maximum gain, i.e., the one which most improves the observation of the “worst-case” state u_K under the expanded observation operator $[L^T, \ell_{K+1}]^T$ and its associated norm. We thereby iteratively approximate the information that would be obtained by measuring with all sensors in the library \mathcal{L} . For fixed θ_K and in combination with selecting u to have the smallest observability in \mathcal{V}_θ , we arrive at an algorithm similar to worst-case orthogonal matching pursuit (cf. [19,20]) but generalized to deal with the covariance function `cov` in the noise model (3).

When expanding L , we also expand the associated noise covariance matrix Σ_L and its Cholesky decomposition. The latter is required for computing the observability coefficient and the observability gain efficiently. The details for this expansion are provided in Algorithm 2 in Section 4.2 below.

Greedy step

We train the observation operator L on all configurations $\theta \in \Xi_{\text{train}}$ by varying for which θ the “worst-case” state is computed. Specifically, we follow a greedy approach where, in iteration K , we identify the configuration θ_K for which the current observation operator L is the least advantageous. To this end, we would ideally choose $\theta_K = \arg \min_{\theta \in \Xi_{\text{train}}} \beta_G(\theta, L)$; however, the computation of $\beta_G(\theta, L)$ for a single θ already requires M full-order model solves, rendering the minimization over Ξ_{train} infeasible. Instead, we therefore approximate $\beta_G(\theta, L) \approx \tilde{\beta}_G(\theta, L)$ using the reduced-order model. The computations are described in Algorithm 4 (`SurrogateObservability`) below. Albeit θ_K might not be the optimal choice in Ξ_{train} for minimizing the full-order $\beta_G(\theta, L)$, with Proposition 2 and the accuracy requirement (17), we are still guaranteed that $\beta_G(\theta_K, L)$ has a similarly small value.

Remark 4. Since the computation of $\tilde{\beta}_G(\theta, L)$ requires as many reduced-order model solves as needed for the posterior covariance matrix over the surrogate model, it is possible to directly target an (approximated) OED utility function in the greedy step in place of $\tilde{\beta}_G(\theta, L)$ without major concessions in the computational efficiency. The data-matching step can then still be performed for the “worst-case” parameter with only one full-order model solve, though its benefit for the utility function should be evaluated carefully.

We proceed to identify the corresponding “worst-case” parameter $\mathbf{m}_K = \arg \min_{\mathbf{m} \in \mathbb{R}^M} \left\| L\tilde{u}_{\theta_K}(\mathbf{m}) \right\|_{\Sigma_L^{-1}} \left\| \mathbf{m} \right\|_{\Sigma_{\text{pr}}}^{-1}$, i.e., the parameter direction for which the least significant observation is achieved. The basis coefficients of \mathbf{m}_K in the eigenvector basis $\{\varphi_m\}_{m=1}^M$ are computed within the call to `SurrogateObservability` (Algorithm 4) with no additional computational effort. Once θ_K has been chosen, \mathbf{m}_K can be assembled in $\mathcal{O}(M^2)$. Similarly to θ_K , \mathbf{m}_K is solely chosen based on the reduced-order surrogate. However, with Proposition 3, the observability $\left\| Lu_{\theta_K}(\mathbf{m}_K) \right\|_{\Sigma_L^{-1}} \left\| \mathbf{m}_K \right\|_{\Sigma_{\text{pr}}}^{-1}$ of \mathbf{m}_K under the full-order model is close to $\beta_G(\theta_K, L)$ such that \mathbf{m}_K may indeed serve as an approximate “worst-case” state for the full-order model.

Termination

Algorithm 1 terminates when $K_{\text{max}} \leq K_{\mathcal{L}}$ sensors have been selected. However, this termination criterion can easily be adapted to prescribe a minimum value of the observability coefficient chosen with respect to the observability $\beta_G(\theta, L)$ achieved with the entire sensor library.

Runtime

Assuming the dominating computational restriction is the model evaluation to solve for $u_\theta(\mathbf{m})$ – as is usually the case for PDE models – then the runtime of each iteration in Algorithm 1 is determined by one full-order model evaluation, and $K_{\mathcal{L}}$ sensor measurements of the full-order state. Compared to computing the posterior covariance matrix for the chosen configuration, the data-matching step saves $M - 1$ full-order model solves.

Algorithm 2: CholeskyExpansion.

Input: observation operator $L = [\ell_1, \dots, \ell_K]^T$, noise covariance matrix Σ_L , Cholesky matrix C_L , new sensor $\ell \in \mathcal{U}'$

$L \leftarrow [\ell_1, \dots, \ell_K, \ell]^T$ // operator expansion

if $K = 0$ **then**

$\Sigma_L \leftarrow (\text{cov}(\ell, \ell)), C_L \leftarrow (\sqrt{\text{cov}(\ell, \ell)}) \in \mathbb{R}^{1 \times 1}$ // first sensor

else

$\mathbf{v} \leftarrow [\text{cov}(\ell_1, \ell), \dots, \text{cov}(\ell_K, \ell)]^T \in \mathbb{R}^K$ // matrix expansion

$\mathbf{w} \leftarrow C_L^{-1} \mathbf{v} \in \mathbb{R}^K, s \leftarrow \text{cov}(\ell, \ell), c \leftarrow s - \mathbf{w}^T \mathbf{w} \in \mathbb{R}$

$\Sigma_L \leftarrow \begin{pmatrix} \Sigma_L & \mathbf{v} \\ \mathbf{v}^T & s \end{pmatrix}, C_L \leftarrow \begin{pmatrix} C_L & \mathbf{0} \\ \mathbf{w}^T & c \end{pmatrix} \in \mathbb{R}^{(K+1) \times (K+1)}$

return L, Σ_L, C_L

The other main factor in the runtime of Algorithm 1 is the $|\Xi_{\text{train}}|M$ reduced-order model evaluations with $K_{\mathcal{L}}$ sensor evaluations each required by the greedy step. Since these are the same in each iteration, they can be pre-computed and stored if memory permits. The parameter dimension M not only enters as a scaling factor, but also affects the cost of the reduced-order model itself since larger values of M generally require larger or more complicated reduced-order models to achieve the desired accuracy (17). In turn, the computational cost of the reduced-order model indicates how large Ξ_{train} may be chosen for a given computational budget. While some cost can be saved through adaptive training sets and models, overall, this connection to M stresses the need for an adequate initial parameter reduction as discussed in Section 3.2.

4.2. Cholesky decomposition

The observability coefficient $\beta_G(\theta, L)$ is connected to the noise model and the covariance function cov through the noise covariance matrix Σ_L whose inverse enters the norm $\|\cdot\|_{\Sigma_L^{-1}}$ and the posterior covariance matrix $\Sigma_{\text{post}}^{L, \theta}$. The inversion poses a challenge when the noise is correlated, i.e., when Σ_L is not diagonal: in this case, even the expansion of L with a single sensor $\ell \in \mathcal{L}$ changes each entry of Σ_L^{-1} . In naive computations of the observability coefficients and the posterior covariance matrix, this leads to M dense linear system solves of order $\mathcal{O}((K+1)^3)$ each time the observation operator is expanded. In the following, we therefore expound on how Σ_L^{-1} changes under expansion of L to exploit its structure when comparing potential sensor choices.

Suppose $L = [\ell_1, \dots, \ell_K]^T$ has already been chosen with sensors $\ell_k \in \mathcal{U}'$, but shall be expanded by another sensor ℓ to

$$[L, \ell] := [\ell_1, \dots, \ell_K, \ell]^T : \mathcal{U} \rightarrow \mathbb{R}^{K+1}.$$

Following definition (4), the noise covariance matrix $\Sigma_{[L, \ell]}$ of the expanded operator $[L, \ell]$ has the form

$$\Sigma_{[L, \ell]} = \begin{pmatrix} \Sigma_L & \mathbf{v}_{L, \ell} \\ \mathbf{v}_{L, \ell}^T & v_{\ell, \ell} \end{pmatrix} = \begin{pmatrix} C_L & \mathbf{0} \\ \mathbf{c}_{L, \ell}^T & c_{\ell, \ell} \end{pmatrix} \begin{pmatrix} C_L^T & \mathbf{c}_{L, \ell} \\ \mathbf{0} & c_{\ell, \ell} \end{pmatrix},$$

where $C_L C_L^T = \Sigma_L \in \mathbb{R}^{K \times K}$ is the Cholesky decomposition of the s.p.d. noise covariance matrix Σ_L for the original observation operator L , and $\mathbf{v}_{L, \ell}, \mathbf{c}_{L, \ell} \in \mathbb{R}^K, v_{\ell, \ell}, c_{\ell, \ell} \in \mathbb{R}$ are defined through

$$\begin{aligned} [\mathbf{v}_{L, \ell}]_i &:= \text{cov}(\ell_i, \ell), & \mathbf{c}_{L, \ell} &:= C_L^{-1} \mathbf{v}_{L, \ell}, \\ v_{\ell, \ell} &:= \text{cov}(\ell, \ell), & c_{\ell, \ell} &:= \sqrt{v_{\ell, \ell} - \mathbf{c}_{L, \ell}^T \mathbf{c}_{L, \ell}}. \end{aligned}$$

Note that $\Sigma_{[L, \ell]}$ is s.p.d. by the assumptions posed on cov in Section 2; consequently, $c_{\ell, \ell}$ is well-defined and strictly positive. With this factorization, the expanded Cholesky matrix $C_{[L, \ell]}$ with $C_{[L, \ell]} C_{[L, \ell]}^T = \Sigma_{[L, \ell]}$ can be computed in $\mathcal{O}(K^2)$, dominated by the linear system solve with the triangular C_L for obtaining $\mathbf{c}_{L, \ell}$. It is summarized in Algorithm 2. We refer to [66], chapter 4, pp. 168-170, for further discussion of the Cholesky decomposition on submatrices.

Using the Cholesky decomposition, the inverse of $\Sigma_{[L, \ell]}$ factorizes to

$$\Sigma_{[L, \ell]}^{-1} = \begin{pmatrix} C_L^T & \mathbf{c}_{L, \ell} \\ \mathbf{0} & c_{\ell, \ell} \end{pmatrix}^{-1} \begin{pmatrix} C_L & \mathbf{0} \\ \mathbf{c}_{L, \ell}^T & c_{\ell, \ell} \end{pmatrix}^{-1} = \begin{pmatrix} C_L^{-T} & \mathbf{r}_{L, \ell} \\ \mathbf{0} & 1/c_{\ell, \ell} \end{pmatrix} \begin{pmatrix} C_L^{-1} & \mathbf{0} \\ \mathbf{r}_{L, \ell}^T & 1/c_{\ell, \ell} \end{pmatrix}, \quad (24)$$

where

$$\mathbf{r}_{L, \ell} := -\frac{1}{c_{\ell, \ell}} C_L^{-T} \mathbf{c}_{L, \ell} = -\frac{1}{c_{\ell, \ell}} C_L^{-T} C_L^{-1} \mathbf{v}_{L, \ell} = -\frac{1}{c_{\ell, \ell}} \Sigma_L^{-1} \mathbf{v}_{L, \ell}.$$

The advantage of the factorization (24) is that the dense Σ_L^{-1} does not need to be explicitly computed when comparing designs, leading to significant computational savings as shown below.

Algorithm 3: ObservabilityGain.

Input: observation operator $L = [\ell_1, \dots, \ell_K]^T$, Cholesky matrix C_L , sensor candidate $\ell \in \mathcal{U}'$, state $u \in \mathcal{U}$

```

d ← Lu, z ← CL-1d                                     // preparation
if K = 0 then
    return ℓ(uK)2/cov(ℓ, ℓ)                             // one sensor only
else
    v ← [cov(ℓ1, ℓ), ..., cov(ℓK, ℓ)]T ∈ ℝK           // general case
    w ← CL-1v ∈ ℝK
    return  $\frac{(\ell(u_K) - w^T z)^2}{\text{cov}(\ell, \ell) - w^T w}$ 

```

4.3. Observability gain

Using (24), for an arbitrary state $u \in \mathcal{U}$, the norm of the extended observation $[L, \ell](u) = [Lu^T, \ell(u)]^T \in \mathbb{R}^{K+1}$ in the corresponding norm $\|\cdot\|_{\Sigma_{[L, \ell]}^{-1}}$ is connected to the original observation $Lu \in \mathbb{R}^K$ in the original norm $\|\cdot\|_{\Sigma_L^{-1}}$ via

$$\begin{aligned}
 \|[L, \ell](u)\|_{\Sigma_{[L, \ell]}^{-1}}^2 &= \begin{pmatrix} Lu \\ \ell(u) \end{pmatrix}^T \begin{pmatrix} \Sigma_L & \mathbf{v}_{L, \ell} \\ \mathbf{v}_{L, \ell}^T & v_{\ell, \ell} \end{pmatrix}^{-1} \begin{pmatrix} Lu \\ \ell(u) \end{pmatrix} \\
 &= \begin{pmatrix} Lu \\ \ell(u) \end{pmatrix}^T \begin{pmatrix} C_L^{-T} & \mathbf{r}_{L, \ell} \\ \mathbf{0} & 1/c_{\ell, \ell} \end{pmatrix} \begin{pmatrix} C_L^{-1} & \mathbf{0} \\ \mathbf{r}_{L, \ell}^T & 1/c_{\ell, \ell} \end{pmatrix} \begin{pmatrix} Lu \\ \ell(u) \end{pmatrix} \\
 &= \begin{pmatrix} C_L^{-1} Lu \\ \mathbf{r}_{L, \ell}^T Lu + \ell(u)/c_{\ell, \ell} \end{pmatrix}^T \begin{pmatrix} C_L^{-1} Lu \\ \mathbf{r}_{L, \ell}^T Lu + \ell(u)/c_{\ell, \ell} \end{pmatrix} \\
 &= (Lu)^T C_L^{-T} C_L^{-1} Lu + (\mathbf{r}_{L, \ell}^T Lu + \ell(u)/c_{\ell, \ell})^2 \\
 &= \|Lu\|_{\Sigma_L^{-1}}^2 + (\mathbf{r}_{L, \ell}^T Lu + \ell_{K+1}(u)/c_{\ell, \ell})^2 \\
 &\geq \|Lu\|_{\Sigma_L^{-1}}^2.
 \end{aligned} \tag{25}$$

We conclude from this result that the norm $\|Lu\|_{\Sigma_L^{-1}}$ of any observation, and therefore also the continuity coefficient γ_L defined in (6), is increasing under expansion of L despite the change in norms. For any configuration θ , the observability coefficient $\beta_G(\theta, L)$ is thus non-decreasing when sensors are selected iteratively. This guarantees in particular, that by iteratively increasing the observability at the “worst-case” parameters and hyper-parameters, we increase the minimum of $\tilde{\beta}_G(\theta, L)$ throughout the training domain.

Given a state $u \in \mathcal{U}$ and an observation operator L , we can determine the sensor $\ell_{K+1} \in \mathcal{L}$ that increases the observation of u the most by comparing the increase $(\mathbf{r}_{L, \ell}^T Lu + \ell(u)/c_{\ell, \ell})^2$ for all $\ell \in \mathcal{L}$. Algorithm 3 summarizes the computation of this observability gain for use in Algorithm 1. Its general runtime is determined by $K+1$ sensor evaluations and two linear solves with the triangular Cholesky matrix C_L in $\mathcal{O}(K^2)$. When called with the same L and the same state u for different candidate sensors ℓ , the preparation step must only be performed once, which reduces the runtime to one sensor evaluation and one linear system solve in all subsequent calls. Compared to computing $\|[L, \ell](u)\|_{\Sigma_{[L, \ell]}^{-1}}^2$ for all $K_{\mathcal{L}}$ candidate sensors in the library \mathcal{L} , we save $\mathcal{O}(K_{\mathcal{L}} K^2)$.

4.4. Computation of the observability coefficient

We next discuss the computation of the observability coefficient $\beta_G(\theta, L)$ for a given configuration θ and observation operator L . Using the eigenvector basis $\{\varphi_m\}_{m=1}^M$ of Σ_{pr} , we define the matrix

$$\mathbf{M}(\theta) := [Lu_\theta(\varphi_1), \dots, Lu_\theta(\varphi_M)] \in \mathbb{R}^{K \times M} \tag{26}$$

featuring all observations of the associated states $u_\theta(\varphi_j)$ for the configuration θ . The observability coefficient $\beta_G(\theta, L)$ can then be computed as the square root of the minimum eigenvalue λ^{\min} of the generalized eigenvalue problem

$$\mathbf{M}(\theta)^T C_L^{-T} C_L^{-1} \mathbf{M}(\theta) \mathbf{m}_{\min} = \lambda^{\min} \mathbf{D}_{\text{pr}}^{-1} \mathbf{m}_{\min}. \tag{27}$$

Note that (27) has M real, non-negative eigenvalues because the matrix on the left is symmetric positive semi-definite, and $\mathbf{D}_{\text{pr}} = \text{diag}(\lambda_{\text{pr}}^1, \dots, \lambda_{\text{pr}}^M)$ is s.p.d. (cf. [66]). The eigenvector \mathbf{m}_{\min} contains the basis coefficients in the eigenvector basis $\{\varphi_m\}_{m=1}^M$ of the “worst-case” parameter, i.e. the infimizer of $\beta_G(\theta, L)$.

The solution of the eigenvalue problem can be computed in $\mathcal{O}(M^3)$, with an additional $\mathcal{O}(MK^2 + M^2K)$ for the computation of the left-hand side matrix in (27). The dominating cost is hidden in $\mathbf{M}(\theta)$ since it requires KM sensor observations and K full-order model solves. To reduce the computational cost, we therefore approximate $\beta_G(\theta, L)$ with $\tilde{\beta}_G(\theta, L)$ by exchanging the full-order states $u_\theta(\varphi_j)$ in (26) with their reduced-order approximations $\tilde{u}_\theta(\varphi_j)$. The procedure is summarized in Algorithm 4.

Algorithm 4: SurrogateObservability.

Input: configuration $\theta \in \mathcal{P}$, observation operator $L = [\ell_1, \dots, \ell_K]^T$ with $K > 0$, Cholesky matrix C_L

$N \leftarrow \min\{M, K\}$ // parameter restriction

$\mathbf{M} \leftarrow [L\tilde{u}_\theta(\varphi_1), \dots, L\tilde{u}_\theta(\varphi_N)], \mathbf{S} \leftarrow [\langle u_\theta(\varphi_i), u_\theta(\varphi_j) \rangle_U]_{i,j=1}^N$ // matrix setup

Find $(\lambda^{\min}, \mathbf{m}_{\min})$ of $[C_L^{-1}\mathbf{M}]^T [C_L^{-1}\mathbf{M}] \mathbf{m}_{\min} = \lambda^{\min} \mathbf{S} \mathbf{m}_{\min}$ // eigenvalue problem

return $\sqrt{\lambda^{\min}}, \mathbf{m}_{\min}$

Remark 5. If $K < M$, Algorithm 4 restricts the parameter space, as discussed in Section 3.2, to the span of the first K eigenvectors $\varphi_1, \dots, \varphi_K$ encoding the least certain directions in the prior. A variation briefly discussed in [19] in the context of the PBDW method to prioritize the least certain parameters even further is to only expand the parameter space once the observability coefficient on the subspace surpasses a predetermined threshold.

Remark 6. For selecting a design L after a reduction in the parameter space (see Section 3.2), we replace the observability coefficient $\beta_G(\theta, L)$ with $\beta_V(\theta, L)$ defined in (16). In this case, we follow the same computational procedure, but exchange the right-hand side matrix $\mathbf{D}_{\text{pr}}^{-1}$ in (27) with the U -inner-product matrix for the states $u_\theta(\varphi_1), \dots, u_\theta(\varphi_M)$ (definitions below). In particular, the result (25) implies that similar to $\beta_G(\theta, L)$, $\beta_V(\theta, L)$ is also non-decreasing under expansion of L , such that the motivations and analysis for the individual steps carry over.

5. Numerical results

We numerically confirm the validity of our sensor selection approach using a geophysical model of a section of the Perth Basin in Western Australia. The basin has raised interest in the geophysics community due to its high potential for geothermal energy, e.g., [67–71]. We focus on a subsection that spans an area of 63 km \times 70 km and reaches 19 km below the surface. The model was introduced in [72] and the presented section of the model was discussed extensively in the context of MOR in [73,74]. In particular, the subsurface temperature distribution is described through a steady-state heat conduction problem with different subdomains for the geological layers, and local measurements may be obtained through boreholes. The borehole locations need to be chosen carefully due to their high costs (typically several million dollars, [75]), which in turn motivates our application of Algorithm 1. For demonstration purposes, we make the following simplifications to our test model: 1) We neglect radiogenic heat production; 2) we merge geological layers with similar conductive behaviors; and 3) we scale the prior to emphasize the influence of different sensor measurements on the posterior. All computations were performed in Python 3.7 on a computer with a 2.3 GHz Quad-Core Intel Core i5 processor and 16 GB of RAM. The source code is available on GitHub (nicolearetz/greedy-sensor-selection).

This section is organized as follows: In Section 5.1, we introduce our modeling assumptions, aiming to provide a comprehensible explanation for our choices. In our first experiment, Section 5.2, we numerically verify whether we can indeed use the observability coefficient to identify sensors that have a similarly small utility function value as the A-, D- and E-optimal choices. Our second experiment, Section 5.3, is designed to show how much the performance of our algorithm is influenced by the sensor library, and should serve as a point of reference for the algorithm’s expected performance in different settings. The purpose of our final experiment, Section 5.4, is to show the scalability of our algorithm to large sensor libraries. In any of the three experiments, we compare two setups of our sensor selection algorithm: In the first, denoted “ θ_{ref} -training” or “proposal, fixed config.”, the algorithm only trains on a reference hyper-parameter θ_{ref} using the full-order model to evaluate how well the algorithm performs in the best-case scenario with no model approximation and no variation in the hyper-parameter. In the second setup, denoted “ Ξ_{train} -training” or “proposal”, the algorithm is trained using a reduced-order surrogate model on a hyper-parameter training set $\Xi_{\text{train}} \not\ni \theta_{\text{ref}}$ just as described in Section 4.

5.1. Model description

The model of this section was originally created to present strategies for a better evaluation of the geothermal potential, which in turn can be used in the decision about potential sites for geothermal power plants. The main uncertainty of the model is the geothermal heat flux of the bottom boundary condition since the model extends far beneath the surface (19 km). Closer to the surface the rock structures influence the temperature so much that the geothermal heat flux exhibits far more local variations, leading to even more uncertainty in the model. The objective for our Bayesian inverse problem is to estimate the geothermal heat flux from temperature measurements and thereby improve the accuracy of the model for better informing and reducing the economical risk of future geothermal power plant placements. The goal for our OED problem is to choose where to take these measurements such that the geothermal heat flux will be approximated well for each combination of conductive properties in the geology.

Hyper-parameterized forward model

We model the subsurface temperature distribution u_θ as our state variable with the steady-state PDE

$$-\nabla \cdot (\theta \nabla u_\theta) = 0 \quad \text{in } \Omega := (0, 0.2714) \times (0, 0.9) \times (0, 1) \subset \mathbb{R}^3. \quad (28)$$

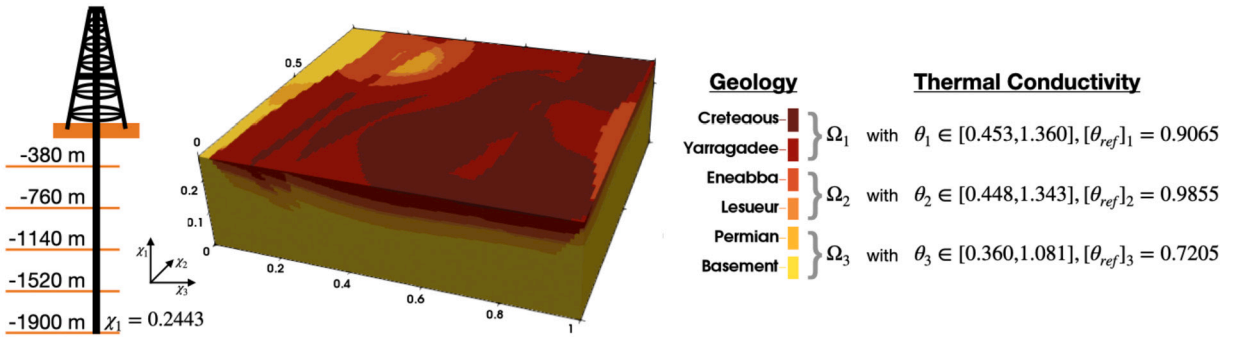


Fig. 1. Schematic overview of the Perth Basin section. Left: drilling depths for potential measurements. The depths were chosen within the first 2 km below surface to reflect the typical depth of hydrocarbon boreholes in the literature [76] while also allowing for point evaluations of the model's subsurface temperature. Note that point evaluations are standard for geophysical models because a borehole (diameter approximately 1 m) is very small compared to the size of the model (in this case 63 km \times 70 km \times 19 km) [77,76]. Middle: geometry Ω with subdomains for the geological layers. Plot adapted from [73]. Right: configuration range and reference values for thermal conductivity θ on each subdomain. The bounds are obtained from the reference values (cf. [72,73]) with a $\pm 50\%$ margin. (For interpretation of the colors in the figure(s), the reader is referred to the web version of this article.)

The domain Ω is a non-dimensionalized representation of the Perth Basin section, and $\theta : \Omega \rightarrow \mathbb{R}_{>0}$ the local thermal conductivity. Geologically, the section has the six geological layers shown in Fig. 1 which have been subjected to several geological processes such as deposition and deformation. For our demonstration purposes here, we group the geological layers further by their main conductive behavior into three subdomains $\Omega_i \subset \Omega$, $\Omega = \bigcup_{i=1,2,3} \Omega_i$. For simulation purposes, the thermal conductivity θ can be considered spatially constant on each layer. In a slight abuse of notation, this lets us identify the field $\theta : \Omega \rightarrow \mathbb{R}_{>0}$ with the vector

$$\theta = (\theta_1, \theta_2, \theta_3) \in \mathcal{P} := [0.453, 1.360] \times [0.448, 1.343] \times [0.360, 1.081],$$

such that $\theta|_{\Omega_i} \equiv \theta_i$. However, while we can consider the position of the layers Ω_i to be fixed as they are often estimated beforehand from, for instance, geophysical campaigns such as seismic surveys, their thermal conductivity θ_i can only be determined within a range, partially because these layers are located deep in the ground. We therefore take the thermal conductivity of each layer as our hyper-parameter, with the hyper-parameter domain $\mathcal{P} \subset \mathbb{R}$ reflecting its variability. The bounds of \mathcal{P} are taken from the literature [72,73] with a $\pm 50\%$ margin.

For the boundary conditions, we impose zero-Dirichlet boundary conditions at the surface, and zero-Neumann ("no-flow") boundary conditions at the lateral faces of the domain. Non-zero Dirichlet boundary conditions obtained from satellite data could be considered via a lifting function and an affine transformation of the measurement data (see [74]). We model the remaining boundary Γ_{In} at the base of the domain as a Neumann boundary condition

$$\mathbf{n} \cdot \nabla u_\theta = g_{\text{flux}} \quad \text{a.e. on } \Gamma_{\text{In}} := \{0\} \times [0, 0.9] \times [0, 1]$$

where $\mathbf{n} : \Gamma_{\text{In}} \rightarrow \mathbb{R}^3$ is the outward pointing unit normal on Ω , and $g_{\text{flux}} : \Gamma_{\text{In}} \rightarrow \mathbb{R}$ the geothermal heat flux.

Modeled uncertainty

The geothermal heat flux at depth can only be observed indirectly and needs to be inferred from, for instance, temperature measurements at boreholes. We therefore treat g_{flux} as uncertain, and impose a parameterization of the form $g_{\text{flux}} = \mathbf{m} \cdot \mathbf{p}$ where $\mathbf{m} \in \mathbb{R}^5$ is a random variable and $\mathbf{p} : (0, 0.9) \times (0, 1) \rightarrow \mathbb{R}$ is a vector composed of quadratic polynomials:

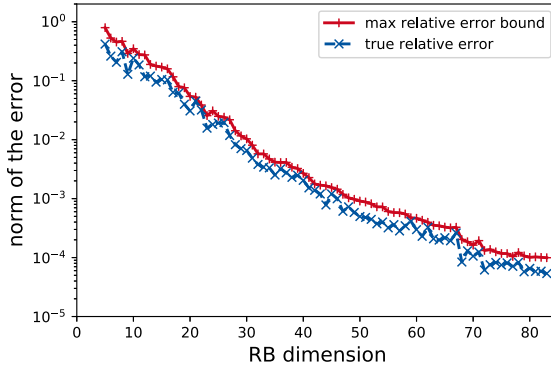
$$\mathbf{p}(x_1, x_2) := \sqrt{\frac{10}{9}} \begin{pmatrix} 1 \\ \sqrt{3}(\frac{20}{9}x_1 - 1) \\ 2x_2 - 1 \\ \sqrt{5}(\frac{200}{27}x_1^2 - \frac{20}{3}x_1 + 1) \\ \sqrt{5}(6x_2^2 - 6x_2 + 1) \end{pmatrix}.$$

The entries of \mathbf{p} have been chosen such that they are orthonormal in the L^2 -norm over $(0, 0.9) \times (0, 1)$. In this parameterization, the geothermal heat flux has a quadratic behavior both in north-south and east-west direction. This setup reflects typical geophysical basal boundary conditions, where it is most common to assume a constant Neumann heat flux (e.g., [73]), and sometimes a linear one (e.g., [72]). With the quadratic functions, we allow an additional degree of freedom than typically considered.

For characterizing the spatial uncertainty in the geothermal heat flux, we model the coefficient vector $\mathbf{m} \in \mathbb{R}^5$ as a random variable $\mathbf{m} \sim \pi_{\text{pr}} = \mathcal{N}(\mathbf{m}_{\text{pr}}, \Sigma_{\text{pr}})$ with

$$\mathbf{m}_{\text{pr}} = [50, 0, 0, 0, 0]^T \in \mathbb{R}^5, \quad \Sigma_{\text{pr}} = \text{diag}(100, 10, 10, 1, 1) \in \mathbb{R}^{5 \times 5}.$$

We have chosen this prior such that the largest uncertainty is attributed to the constant entry in \mathbf{p} , and the quadratic terms are treated as the most certain with prior zero. This choice reflects that the basal boundary is so deep – 19 km below the surface – that



Reduced-order model

RB dimension	83
training time	37.58 min
training accuracy	1e-4
RB solve	0.97 ms
↪ speedup	3,058
RB error bound	4.78 ms
↪ speedup	515

Fig. 2. Training of the RB surrogate model for the Perth Basin section using a greedy algorithm, cf. [61]. In each iteration, the algorithm samples θ randomly, computes an error bound in the form (29), and then expands the RB space with the full-order solution for which the error bound was largest. The training was concluded when the target accuracy (29) was reached for 511,000 consecutively drawn hyper-parameters. The RB model is used below for the greedy hyper-parameter selection within Algorithm 1. Left: Maximum relative error bound (29) in the course of the greedy algorithm and corresponding true relative error at the configuration θ chosen for space expansion. On the right: Performance pointers for the obtained RB model after the target accuracy (29) was reached; online computation times and speedups are averages computed over 1000 randomly drawn configurations θ .

local variations in the geothermal heat flux have mostly stabilized. Our goal is to choose sensor locations in the form of borehole positions and drilling depths to reliably estimate \mathbf{m} and consequently the geothermal heat flux g_{flux} for any admissible realization of the thermal conductivity θ .

Discretization

The problem is discretized using a finite element (FE) space \mathcal{U} of dimension 132,651 with piece-wise linear basis functions. The underlying mesh was created with GemPy ([78]) and MOOSE ([79]). We equip \mathcal{U} with the inner product $\langle u, \phi \rangle_{\mathcal{U}} := \int_{\Omega} \nabla u \cdot \nabla \phi d\Omega$. Note that $\langle \cdot, \cdot \rangle_{\mathcal{U}}$ is indeed an inner product due to the Dirichlet boundary conditions. Owing to the structure of the governing equation (28) and the division of the domain $\Omega = \bigcup_{i=1,2,3} \Omega_i$ into disjoint subdomains, the FE matrices decouple in θ ; we therefore precompute and store an affine decomposition using the library DwarfElephant ([73]) from which any θ -dependent FE stiffness matrix can then be reconstructed. Given a configuration θ and a coefficient vector \mathbf{m} for the geothermal heat flux at Γ_{In} , the FE matrix assembly and the computation of a full-order solution $u_{\theta}(\mathbf{m}) \in \mathcal{U}$ then takes 2.96 s on average.

Reduced-order model

An important requirement of the sensor selection Algorithm 1 is the availability of a reduced-order surrogate model: It enables the algorithm to deal efficiently with the computational burden of comparing designs for various hyper-parameters despite the computationally expensive PDE model. Here, we chose a reduced basis (RB) method to exploit the model's affine decomposition further.

The RB model was trained with a greedy algorithm (cf. [61,80]): Using a rigorous and certified *a posteriori* error bound $\Delta(\theta)$ (cf. [56–60]), we prescribe the relative target accuracy

$$\max_{\mathbf{m} \in \mathbb{R}^M} \frac{\|u_{\theta}(\mathbf{m}) - \tilde{u}_{\theta}(\mathbf{m})\|_{\mathcal{U}}}{\|u_{\theta}(\mathbf{m})\|_{\mathcal{U}}} \leq \max_{\mathbf{m} \in \mathbb{R}^M} \frac{\Delta(\theta)}{\|\tilde{u}_{\theta}(\mathbf{m})\|_{\mathcal{U}}} < \varepsilon := 10^{-4} \quad (29)$$

to be reached for 511,000 consecutively drawn, uniformly distributed samples of θ . The training phase and final computational performance of the RB surrogate model are summarized in Fig. 2. The speedup of the surrogate model (approximately a factor of 3,000 without error bounds) justifies its offline training time, with computational savings expected already after 152 approximations of $\beta_G(\theta, L)$.

Temperature measurements

In practice, measurements in typical geothermal data sets are often made in boreholes drilled for hydrocarbon exploration. Here, the maximum temperature is measured but then associated to the bottom of the borehole, known as a “bottom hole temperature measurement”. Therefore the measurements are considered as low quality data subjected to uncertainties [81]. We therefore model a sensor measurement as a point evaluation of the subsurface temperature at the bottom of any admissible borehole (defined below). For sensor selection, we then pose the additional combinatorial restriction that no drilling site may be chosen twice, i.e., boreholes may not overlap.

Sensor libraries

Since boreholes are very expensive [75,82], it is unrealistic that the position and drilling depth of the boreholes can be chosen entirely for the purposes of inferring the geothermal heat flux. More realistically, an expert committee would decide on a candidate set of possible drilling locations first, and the OED for the inference of the geothermal heat flux might be used as a tiebreaker, to decide on the drilling order, or to otherwise inform the opinion of the scientific expert on the committee. For our demonstration

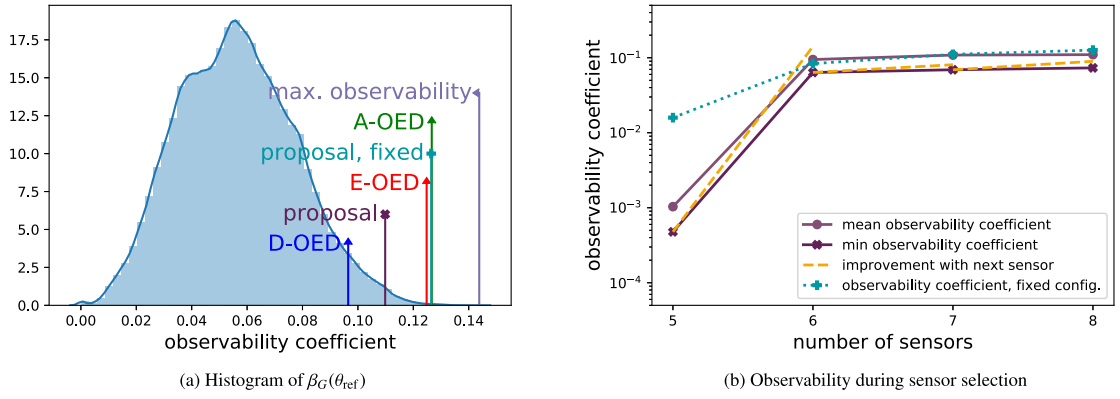


Fig. 3. Observability coefficient when choosing 8 out of the 25 sensor locations in library $\mathcal{L}_{5 \times 5}$. Left: Distribution of the observability coefficient $\beta_G(\theta_{\text{ref}}, L)$ at the reference configuration θ_{ref} over all possible sensor combinations in L . Indicators show the observability coefficients for the A-, D-, and E-optimal choices, the sensor combination with maximum observability, and the sensors chosen by the Algorithm 1 with Ξ_{train} -training (“proposal”, purple, marked “x”) and θ_{ref} -training (“proposal, fixed”, turquoise, marked “+”). Note that the height of the indicator line was chosen solely for readability. Right: Performance of Algorithm 1 over the number of selected sensors. Continuous lines show the minimum (marked “x”) and mean (marked “o”) of $\hat{\beta}_G(\theta, L)$ over Ξ_{train} when training with variable hyper-parameters (Ξ_{train} -training). For each iteration, the dashed lines show the gain in observability achieved with the next sensor for the worst-case configuration. For comparison, the dotted line shows $\beta_G(\theta_{\text{ref}}, L)$ when running Algorithm 1 using the full-order model on the reference configuration θ_{ref} only (θ_{ref} -training).

purposes here, the expert committee may choose between 2,209 potential drilling sites, located on a 47×47 grid over the surface. At each location, a single borehole may be drilled which may reach any one of five depths (see Fig. 1), resulting in a total of 11,045 admissible sensor choices \mathcal{L}_{all} prior to the committee’s selection. For each of our three numerical experiments, we mimic the committee’s choice and choose a subset of $\mathcal{L} \subset \mathcal{L}_{\text{all}}$ as our available sensor library in order to evaluate the performance of Algorithm 1 in different settings.

Noise model

We model the noise covariance between sensor measurements $\ell_x, \ell_{\bar{x}} \in \mathcal{L}_{\text{all}}$ at points $x, \bar{x} \in \Omega$ via

$$\text{cov}(\ell_x, \ell_{\bar{x}}) := a + b - y(h)$$

with the exponential variogram model

$$y(h) := a + (b - a) \left(\frac{3}{2} \max\left\{\frac{h}{c}, 1\right\} - \frac{1}{2} \max\left\{\frac{h}{c}, 1\right\}^3 \right)$$

where $h^2 := (x_2 - \bar{x}_2)^2 + (x_3 - \bar{x}_3)^2$ is the horizontal distance between the points and

$$\begin{aligned} a &:= 2.2054073480730403 && \text{(sill)} \\ b &:= 1.6850672040263555 && \text{(nugget)} \\ c &:= 20.606782733391228 && \text{(range)} \end{aligned}$$

The covariance function was computed via kriging (cf. [83]) from the existing measurements [84]. With this covariance function, the noise between measurements at any two sensor locations is increasingly correlated the closer they are on the horizontal plane. Note that for any subset of sensor locations, the associated noise covariance matrix remains regular as long as each sensor is placed at a distinct drilling location.

5.2. Restricted library

To test the feasibility of the observability coefficient for sensor selection, we first consider a small sensor library (denoted as $\mathcal{L}_{5 \times 5}$ below) with 25 drilling locations positioned on a 5×5 grid. We consider the problem of choosing 8 pair-wise different, unordered sensor locations out of the given 25 positions; this is a combinatorial problem with 1,081,575 possible combinations.

Sensor selection

We run Algorithm 1, using the RB surrogate model and a training set $\Xi_{\text{train}} \subset \mathcal{P}$ with 512,000 configurations on an $80 \times 80 \times 80$ regular grid on \mathcal{P} . When new sensors are chosen, the surrogate observability coefficient $\hat{\beta}_G(\theta, L)$ increases monotonously with a strong incline just after the initial $M = 5$ sensors, followed by a visible stagnation (see Fig. 3b) as is often observed for similar OMP-based sensor selection algorithms (e.g., [19,85,86,18]). Algorithm 1 terminates in 7.93 min with a minimum reduced-order observability of $\hat{\beta}_G(\theta, L) = 0.073227$ and an average of 0.10995. At the reference configuration θ_{ref} , the full-order observability coefficient is $\beta_G(\theta_{\text{ref}}, L) = 1.0985$, slightly below the reduced-order average. We call this training procedure “ Ξ_{train} -training” hereafter and denote the chosen sensors as “ Ξ_{train} -trained sensor set” in the subsequent text and as “proposal” in the plots.

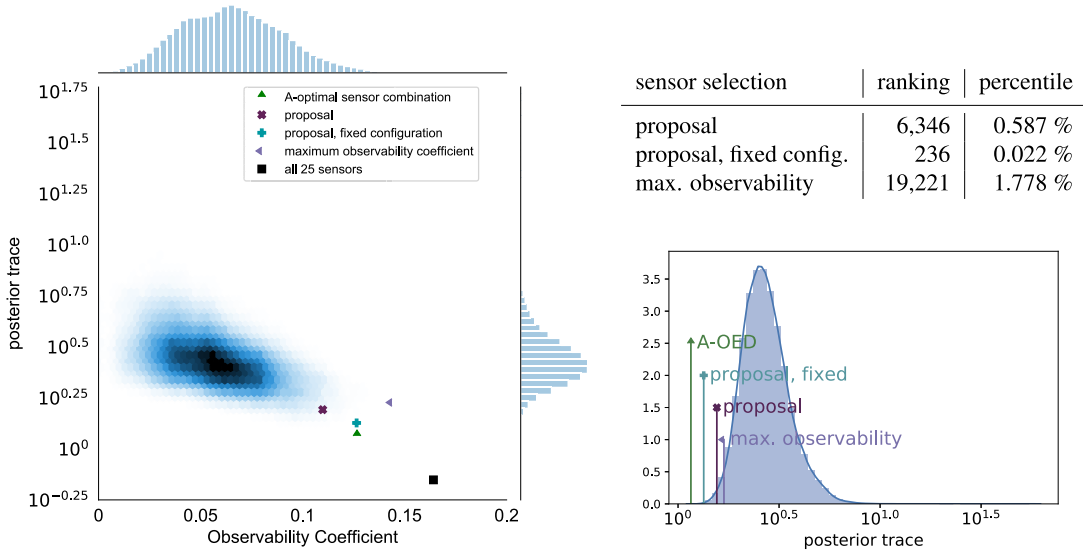


Fig. 4. Joint distribution of $\text{trace}(\Sigma_{\text{post}}^{L,\theta})$ and $\beta_G(\theta_{\text{ref}}, L)$ for $\theta = \theta_{\text{ref}}$ over all 1,081,575 combinations for choosing L with 8 out of the 25 sensor locations in $\mathcal{L}_{5 \times 5}$. For reference, the design $L = \mathcal{L}_{5 \times 5}$ containing all 25 sensors is marked. The marginal distribution for $\text{trace}(\Sigma_{\text{post}}^{L,\theta})$ (vertical axis) is provided on the right for closer inspection. The marginal distribution for $\beta_G(\theta_{\text{ref}}, L)$ (horizontal axis) is provided in Fig. 3a. Indicators highlight the extremal values $(\beta_G(\theta_{\text{ref}}, L), \text{trace}(\Sigma_{\text{post}}^{L,\theta}))$ obtained for the A-optimal design L , and the design with maximal observability. The indicator “proposal, fixed configuration” shows the values $(\beta_G(\theta_{\text{ref}}, L), \text{trace}(\Sigma_{\text{post}}^{L,\theta}))$ when L is obtained with θ_{ref} -training using Algorithm 1. For the indicator “proposal”, L was chosen with Ξ_{train} -training using the RB surrogate model. Out of the 1,081,575 possible choices for L , only 236 designs (0.022%) have an equal or smaller value for $\text{trace}(\Sigma_{\text{post}}^{L,\theta})$ than obtained with θ_{ref} -training. For Ξ_{train} -training, this number increases to 6,346 designs (0.587%), and to 19,221 designs (1.778%) for the design achieving the maximum observability coefficient $\beta_G(\theta_{\text{ref}}, L)$. These numbers are provided in the table (top right) as reference for how well Algorithm 1 and the observability coefficient perform for selecting sensors that are close to being A-optimal.

In order to get an accurate understanding of how the surrogate model $\tilde{u}_\theta(\mathbf{m})$ and the large configuration training set Ξ_{train} influence the sensor selection, we run Algorithm 1 again, this time restricted on the full-order FE model $u_{\theta_{\text{ref}}}(\mathbf{m})$ at only the reference configuration θ_{ref} . The increase in $\beta_G(\theta_{\text{ref}}, L)$ in the course of the algorithm is shown in Fig. 3b. The curve starts significantly above the average for Ξ_{train} -training, presumably because conflicting configurations cannot occur, e.g., when one sensor would significantly increase the observability at one configuration but cause little change in another. However, in the stagnation phase, the curve comes closer to the average achieved with Ξ_{train} -training. The computation finishes within 12.53 s, showing that the long runtime before can be attributed to the size of Ξ_{train} . The final observability coefficient with 8 sensors is $\beta_G(\theta_{\text{ref}}, L) = 0.12647$, above the average over $\tilde{\beta}_G(\theta, L)$ achieved training on Ξ_{train} . We call this training procedure “ θ_{ref} -training” hereafter, and the sensor configuration “ θ_{ref} -trained” in the text or “proposal, fixed config.” in the plots.

Comparison at the reference configuration

For comparing the performance of the Ξ_{train} - and θ_{ref} -trained sensor combinations, we compute – at the reference configuration θ_{ref} – all 1,081,575 posterior covariance matrices $\Sigma_{\text{post}}^{\theta_{\text{ref}}, L}$ for all unordered combinations L of 8 distinct sensors in the sensor library $\mathcal{L}_{5 \times 5}$. For each matrix, we compute the trace (A-OED criterion), the determinant (D-OED criterion), the maximum eigenvalue (E-OED criterion), and the observability coefficient $\beta_G(\theta_{\text{ref}}, L)$. This lets us identify the A-, D-, and E-optimal sensor combinations. The total runtime for these computations is 4 min – well above the 12.53 s of θ_{ref} -training. The (almost) 8 min for Ξ_{train} -training remain reasonable considering it is trained on $|\Xi_{\text{train}}| = 512,000$ configurations and not only θ_{ref} .

A histogram for the distribution of $\beta_G(\theta_{\text{ref}}, L)$ is given in Fig. 3a with markers for the values of the A-, D-, and E-optimal choices and the Ξ_{train} - and θ_{ref} -trained observation operators. Out of these five, the D-optimal choice has the smallest value, since the posterior determinant is influenced less by the maximum posterior eigenvalue and hence the observability coefficient. In contrast, both the A- and E-optimal sensor choices are among the 700 combinations with the largest $\beta_G(\theta_{\text{ref}}, L)$ (this corresponds to the top 0.065%). The θ_{ref} -trained sensors have similar observability and are even among the top 500 combinations. For the Ξ_{train} -trained sensors, the observability coefficient is smaller, presumably because Ξ_{train} -training is not as optimized for θ_{ref} . Still, it ranks among the top 0.705% of sensor combinations with the largest observability.

In order to visualize the connection between the observability coefficient $\beta_G(\theta_{\text{ref}}, L)$ and the classic A-, D-, and E-OED criteria, we plot the distribution of the posterior covariance matrix’s trace, determinant, and maximum eigenvalue over all sensor combinations against $\beta_G(\theta, L)$ in Figs. 4, 5, 6. Overall we observe a strong correlation between the respective OED criteria and $\beta_G(\theta_{\text{ref}}, L)$: It is the most pronounced in Fig. 6 for E-optimality, and the least pronounced for D-optimality in Fig. 5. For all OED criteria, the correlation becomes stronger for smaller scaling factors σ^2 and weakens for large σ^2 when the prior is prioritized (plots not shown). This behavior aligns with the discussion in Section 3.1 that $\beta_G(\theta, L)$ primarily targets the largest posterior eigenvalue and is most decisive for priors with higher uncertainty.

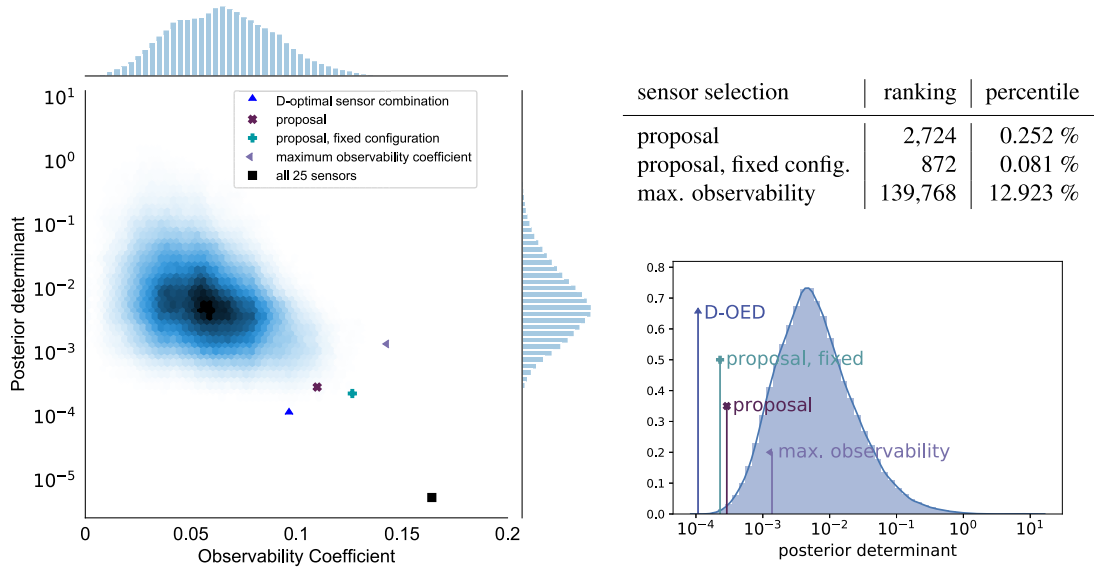


Fig. 5. Joint distribution of the posterior determinant $\det(\Sigma_{\text{post}}^{L,\theta})$ for $\theta = \theta_{\text{ref}}$. See Fig. 4 for details about the plot structure.

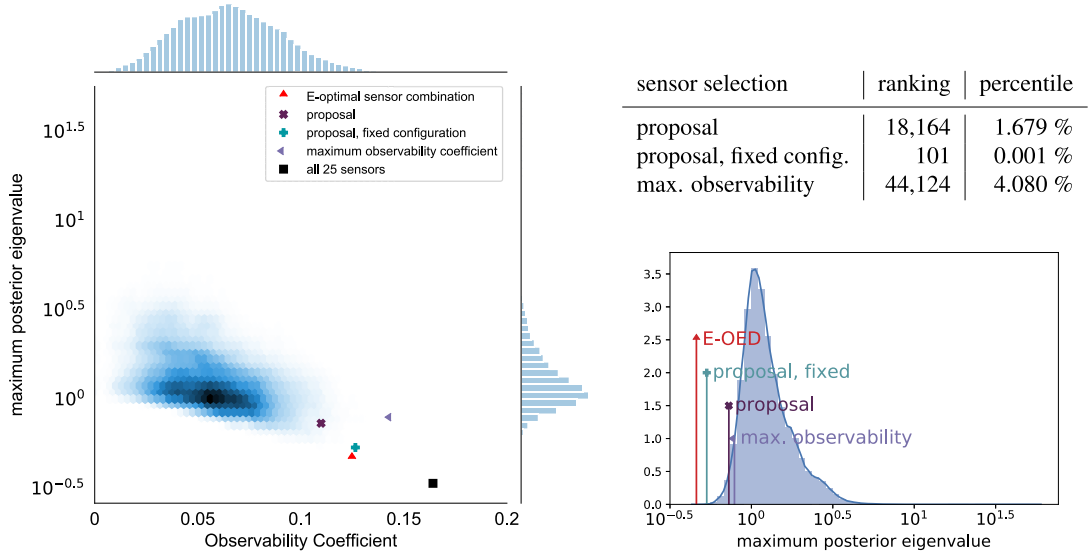
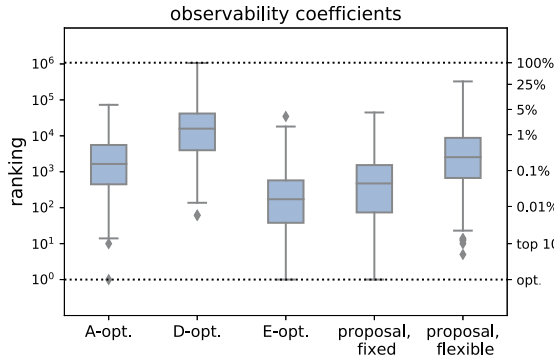


Fig. 6. Joint distribution of the maximum eigenvalue of the posterior covariance matrix $\Sigma_{\text{post}}^{L,\theta}$ for $\theta = \theta_{\text{ref}}$. See Fig. 4 for details about the plot structure. Note that the θ_{ref} -trained sensor combination has the 101-st smallest maximum posterior eigenvalue among all 1,081,575 possibilities.

5.3. Comparison for different libraries

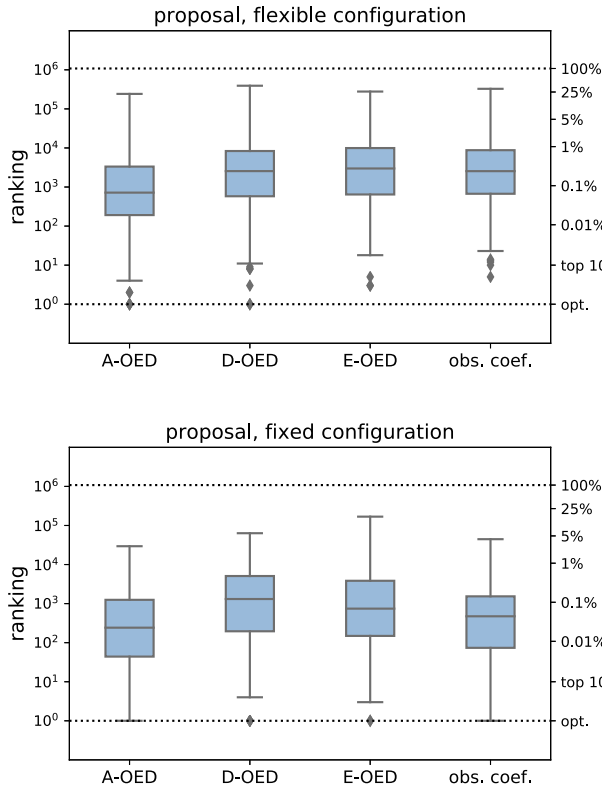
We evaluate the influence of the library $\mathcal{L}_{5 \times 5}$ on our results. To this end, we randomly select 200 sets of new measurement positions from \mathcal{L}_{all} , each consisting of 25 drilling locations with an associated drilling depth. For each library, we run Algorithm 1 to choose 8 sensors, once with Ξ_{train} -training on the surrogate model, and once with the full-order model at θ_{ref} only. For comparison, we then consider in each library each possible combination of choosing 8 unordered sensor sets and compute the trace, determinant, and maximum eigenvalue of the associated posterior covariance matrix at the reference configuration θ_{ref} together with its observability coefficient. This lets us identify the A-, D-, and E-optimal sensor combinations.

Fig. 7 shows how $\beta_G(\theta_{\text{ref}}, L)$ is distributed over the 200 libraries, with percentiles provided in the adjacent table. For 75% of the libraries, the A- and E-optimal, and the Ξ_{train} - and θ_{ref} -trained sensor choices rank among the top 1% of combinations with the largest observability. Due to its non-optimized training for θ_{ref} , the Ξ_{train} -trained sensor set performs slightly worse than what is achieved with θ_{ref} -training, but still yields a comparatively large value for $\beta_G(\theta_{\text{ref}}, L)$. In contrast, overall, the D-optimal sensor choices have smaller observability coefficients, presumably because the minimization of the posterior determinant is influenced less by the maximum posterior eigenvalue.



pctl	design criterion			training	
	A-OED	D-OED	E-OED	θ_{ref}	Ξ_{train}
99-th	3.5835	81.8508	1.1724	2.2223	9.2512
95-th	2.2747	26.8430	0.3601	0.7846	4.0374
75-th	0.5141	3.8600	0.0532	0.1419	0.8106
50-th	0.1527	1.4641	0.0159	0.0438	0.2354
25-th	0.0414	0.3669	0.0035	0.0068	0.0621

Fig. 7. Ranking of $\beta_G(\theta_{\text{ref}}, L)$ for the A-, D-, E-optimal and the θ_{ref} - and Ξ_{train} -trained sensor choices. For any library, the ranking is computed by comparing $\beta_G(\theta_{\text{ref}}, L)$ for all possible sensor combinations. For this experiment, each library contained 25 sensors from which 8 unordered sensors should be chosen, resulting in 1,081,575 possibilities. The distributions (left) were obtained over 200 random sensor libraries. The table, (right), shows the worst-case ranking (in percent) of the corresponding percentiles (“pctl”). For example, for 95% of the 200 random libraries the observability coefficient $\beta_G(\theta_{\text{ref}}, L)$ of the E-optimal experimental design L was among the largest 0.3601% of all designs, indicating that maximizing $\beta_G(\theta, L)$ might indeed be favorable for E-optimal designs. In contrast, 95% of the D-optimal designs ranked only among the largest 26.843% observability coefficients, supporting the claim that maximizing $\beta_G(\theta, L)$ is less suitable for distinguishing between sensors when the goal is D-optimality.



pctl	A-OED	D-OED	E-OED	$\beta_G(\theta_{\text{ref}}, L)$
99-th	3.9240	6.2372	10.8391	9.2512
95-th	1.9093	3.1544	4.5583	4.0374
75-th	0.3083	0.7718	0.9185	0.8106
50-th	0.0664	0.2361	0.2763	0.2354
25-th	0.0177	0.0536	0.0596	0.0621

pctl	A-OED	D-OED	E-OED	$\beta_G(\theta_{\text{ref}}, L)$
99-th	2.5261	2.9752	11.1534	2.2223
95-th	1.0134	1.8324	2.8458	0.7846
75-th	0.1155	0.4698	0.3549	0.1419
50-th	0.0224	0.1212	0.0687	0.0438
25-th	0.0041	0.0181	0.0138	0.0068

Fig. 8. Ranking of the posterior covariance matrix $\Sigma_{\text{post}}^{\theta_{\text{ref}}, L}$ in terms of the A-, D-, E-OED criteria and the observability coefficient $\beta_G(\theta_{\text{ref}}, L)$ when the observation operator $G_{L, \theta}$ is chosen with Algorithm 1 and Ξ_{train} -training (top) or θ_{ref} -training (bottom). The ranking is obtained by comparing all possible unordered combinations of 8 sensors in each sensor library. On the left: Boxplots of the ranking over 200 random sensor libraries; on the right: worst-case ranking (in percent) among different percentiles. Example: For 95% of the libraries, the sensors chosen by Algorithm 1 and θ_{ref} -training were among the 2.88458% of designs with smallest maximum posterior eigenvalue at the reference configuration θ_{ref} (E-OED utility criterion, bottom table). For Ξ_{train} -training, the sensors are not optimized for θ_{ref} specifically, but 95% of the chosen designs were still among the 4.5583% of combinations with smallest posterior eigenvalue.

The ranking of the Ξ_{train} - and θ_{ref} -trained sensor configurations in terms of the posterior covariance matrix’s trace, determinant, and maximum eigenvalue over the 200 libraries is given in Fig. 8. Both perform well and lie for 75% of the libraries within the top 1% of combinations. As the ranking is performed for the configuration parameter θ_{ref} , the θ_{ref} -trained sensor combination performs better, remaining in 95% of the libraries within the top 5% of sensor combinations.

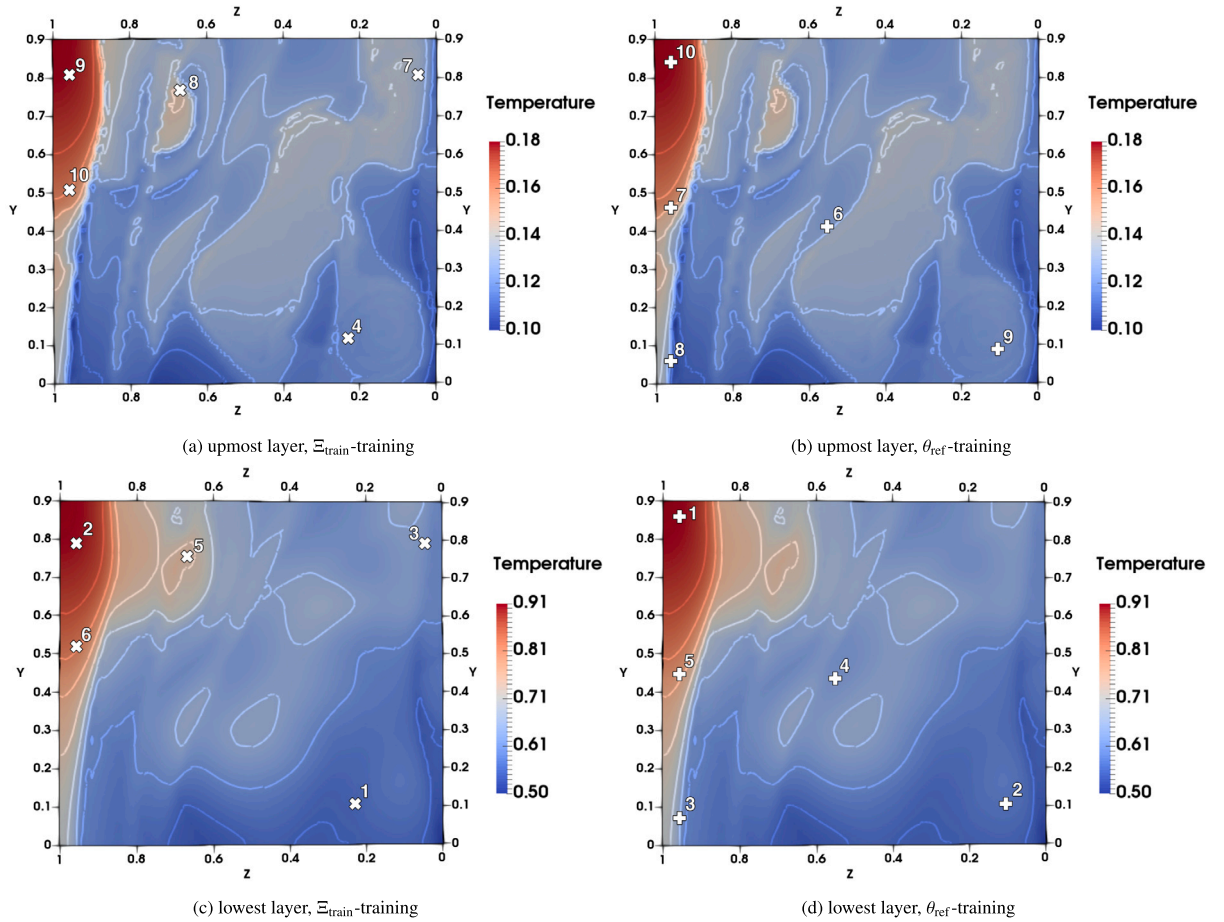


Fig. 9. Sensor positions chosen by Algorithm 1 from a grid of 47×47 available horizontal positions with available 5 depths each, though only the lowest (bottom row) and upmost (top row) layers were chosen. The underlying plot shows cuts through the full-order solution $u_\theta(\mathbf{m})$ at $\theta = \theta_{\text{ref}}$. Top: upmost layer at depth 380 m. Bottom: lowest layer available for measurements at 1.9 km below surface. Left: Sensor positions chosen with Ξ_{train} -training using the RB surrogate model on a training set $\Xi_{\text{train}} \subset \mathcal{P}$ with 10,000 random configurations; runtime 14.19 s for 10 sensors, excluding training of the RB surrogate model (Fig. 2). Right: Sensor positions chosen with θ_{ref} -training using the full-order model at the reference parameter θ_{ref} ; runtime 15.85 s for 10 sensors, including full-order model solves.

5.4. Unrestricted library

We next verify the scalability of Algorithm 1 to large sensor libraries by permitting all 2,209 drilling locations, at each of which at most one measurement may be taken at any of the 5 available measurement depths. Choosing 10 unordered sensors yields approximately 7.29×10^{33} possible combinations. Using the RB surrogate model from before, we run Algorithm 1 once on a training grid $\Xi_{\text{train}} \subset \mathcal{P}$ consisting of 10,000 randomly chosen configurations using only the surrogate model (runtime 14.19 s), and once on the reference configuration θ_{ref} using the full-order model (runtime 15.85 s) for comparison. We terminate the algorithm whenever 10 sensors are selected. Compared to the training time on $\mathcal{L}_{5 \times 5}$ before, the results confirm that the size of the library itself has little influence on the overall runtime but that the full-order computations and the size of Ξ_{train} relative to the surrogate compute dominate.

The sensors chosen by the two runs of Algorithm 1 are shown in Fig. 9. They share many structural similarities:

- **Depth:** Despite the availability of 5 measurement depths, sensors have only been chosen on the lowest and the upmost layers with 5 sensors each. The lower sensors were chosen first (with one exception, sensor 3 in θ_{ref} -training), presumably because the lower layer is closer to the uncertain Neumann boundary condition and therefore yields larger measurement values.
- **Pairing** Each sensor on the lowest layer has a counterpart on the upmost layer that has almost the same position on the horizontal plane. This pairing targets noise sensitivity: With the prescribed error covariance function, the noise in two measurements is increasingly correlated the closer the measurements lie horizontally, independent of their depth coordinate. Choosing a reference measurement near the zero-Dirichlet boundary at the surface helps filter out noise terms in the lower measurement.
- **Organization** On each layer, the sensors are spread out evenly and approximately aligned in 3 rows and 3 columns. The alignment helps distinguish between the constant, linear, and quadratic parts of the uncertain Neumann flux function in north-south and east-west directions.

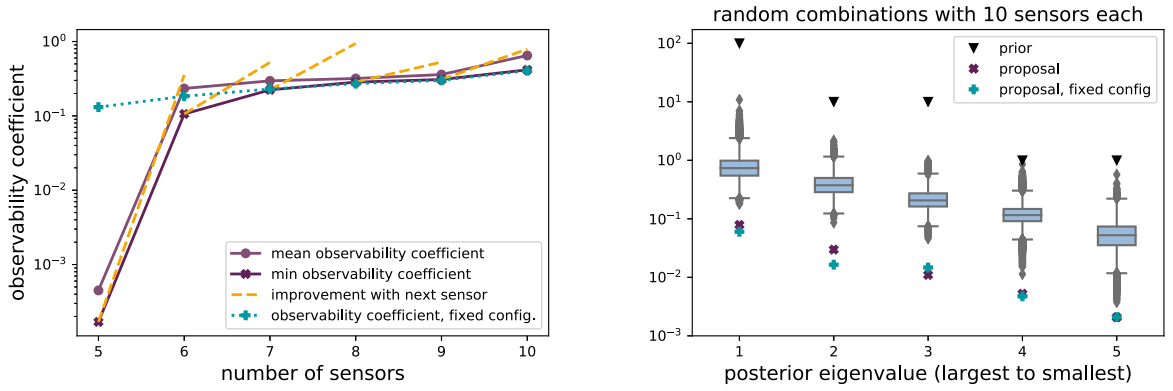


Fig. 10. Left: Observability coefficients during sensor selection with Ξ_{train} - and θ_{ref} -training for a library with 11,045 measurement positions and combinatorial restrictions. Shown are 1) the minimum and mean surrogate observability coefficient $\tilde{\beta}_G(\theta, L)$ over a training set with 10,000 random configurations (Ξ_{train} -training) with final values $\min_{\theta} \tilde{\beta}_G(\theta, L) = 0.4160$ and $\text{mean}_{\theta} \tilde{\beta}_G(\theta, L) = 0.6488$, and 2) the full-order observability coefficient $\beta_G(\theta_{\text{ref}}, L)$ when training on the reference parameter θ_{ref} alone θ_{ref} -training with final value $\beta_G(\theta_{\text{ref}}, L) = 0.4042$. The dashed lines show the improvement achieved by the new sensor at the worst-case configuration during Ξ_{train} -training. Right: Boxplots for the 5 eigenvalues of the posterior covariance matrix $\Sigma^{L,\theta}_{\text{post}}$ at $\theta = \theta_{\text{ref}}$ over 50,000 sets of 10 random sensors. The sensors were chosen uniformly from a $5 \times 47 \times 47$ grid with imposed combinatorial restrictions. The eigenvalues are compared according to their order from largest to smallest. Indicated are also the eigenvalues for the Ξ_{train} -trained (purple, “x”-marker) and θ_{ref} -trained (turquoise, “+”-marker) sensors from Fig. 9. The comparison with the eigenvalues of the prior covariance matrix Σ_{pr} (black, triangular marker), shows the reduction in uncertainty achieved by the different designs.

Fig. 10 (left side) shows the increase in the observability coefficients $\tilde{\beta}_G(\theta, L)$ (for Ξ_{train} -training) and $\beta_G(\theta_{\text{ref}}, L)$ (for θ_{ref} -training) over the number of chosen sensors. We again observe a strong initial incline followed by stagnation for the Ξ_{train} -trained sensors, whereas the curve for θ_{ref} -training already starts at a large value to remain then almost constant. The latter is explained by the positions of the first 5 sensors in Fig. 9 (right), as they are already spaced apart in both directions for the identification of quadratic polynomials. In contrast, for Ξ_{train} -training, the “3 rows, 3 columns” structure is only completed after the sixth sensor (cf. Fig. 9, left). With 6 sensors, the observability coefficients in both training schemes have already surpassed the final observability coefficients with 8 sensors in the previous training on the smaller library $\mathcal{L}_{5 \times 5}$. The final observability coefficients at the reference parameter θ_{ref} are $\beta_G(\theta_{\text{ref}}, L) = 0.4042$ for θ_{ref} -training, and $\beta_G(\theta_{\text{ref}}, L) = 0.3595$ for Ξ_{train} -training.

As a final experiment, we compare the eigenvalues of the posterior covariance matrix $\Sigma^{L,\theta_{\text{ref}}}_{\text{post}}$ for the Ξ_{train} - and θ_{ref} -trained sensors against 50,000 sets of 10 random sensors each. We confirm that all 50,000 sensor combinations comply with the combinatorial restrictions. Boxplots of the eigenvalues are provided in Fig. 10 (right side). We compare the largest eigenvalue of one matrix to the largest eigenvalue of another, the second largest to the second largest, and so on. The eigenvalues of the posterior covariance matrix with sensors chosen by Algorithm 1 are clearly smaller than all posterior eigenvalues for the random sensor combinations, on average by at least factor 10 for each eigenvalue. The comparison with the eigenvalues of the prior covariance matrix also shows how much the uncertainty has been reduced in total.

6. Conclusion

In this work, we analyzed the connection between the observation operator and the eigenvalues of the posterior covariance matrix in the inference of an uncertain parameter via Bayesian inversion for a linear, hyper-parameterized forward model. We identified an observability coefficient whose maximization decreases the uncertainty in the posterior probability distribution for all hyper-parameters. To this end, we proposed a sensor selection algorithm that expands an observation operator iteratively to guarantee a uniformly large observability coefficient for all hyper-parameters. Computational feasibility is retained through a reduced-order model in the greedy step and a data-matching step for the next sensor that only requires a single full-order model evaluation. The validity of the approach was demonstrated on a large-scale heat conduction problem over a section of the Perth Basin in Western Australia. Future extensions of this work are planned to address 1) high-dimensional parameter spaces through parameter reduction techniques, 2) the combination with the PBDW *inf-sup*-criterion to inform sensors by functional analytic means in addition to the noise covariance, and 3) the expansion to non-linear models through a Laplace approximation.

CRediT authorship contribution statement

Nicole Aretz: Conceptualization, Formal analysis, Investigation, Methodology, Software, Validation, Visualization, Writing – original draft, Writing – review & editing. **Peng Chen:** Methodology, Writing – review & editing. **Denise Degen:** Methodology, Software, Visualization, Writing – review & editing. **Karen Veroy:** Conceptualization, Funding acquisition, Methodology, Project administration, Supervision, Writing – review & editing.

Declaration of competing interest

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Co-author serves in an editorial capacity for the journal (K. V.).

Data availability

The source code for the proposed algorithm is available in a public GitHub repository (nicolearetz/greedy-sensor-selection). The large-scale Perth Basin model is available upon request.

Acknowledgements

We would like to thank Tan Bui-Thanh, Youssef Marzouk, Francesco Silva, Andrew Stuart, Dariusz Ucinski, and Keyi Wu for very helpful discussions, and Florian Wellmann at the Institute for Computational Geoscience, Geothermics and Reservoir Geophysics at RWTH Aachen University for providing the Perth Basin Model. This work was supported by the Excellence Initiative of the German federal and state governments and the German Research Foundation through Grants GSC 111 and 33849990/GRK2379 (IRTG Modern Inverse Problems). This project has also received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement n° 818473), the US Department of Energy (grants DE-SC0021239, DE-SC002317), and the US Air Force Office of Scientific Research (grant FA9550-21-1-0084). Peng Chen is partially supported by the NSF grants #2245674, #2245111, and #2325631.

References

- [1] A.M. Stuart, Inverse problems: a Bayesian perspective, *Acta Numer.* 19 (2010) 451–559.
- [2] T. Bui-Thanh, O. Ghattas, Bayes is optimal, *ICES Report* 15, 2015.
- [3] A. Zellner, Optimal information processing and Bayes's theorem, *Am. Stat.* 42 (1988) 278–280.
- [4] I. Stober, K. Bucher Geothermie, Springer, 2012.
- [5] D. Ucinski, Optimal Measurement Methods for Distributed Parameter System Identification, CRC Press, 2004.
- [6] V.B. Melas, Functional Approach to Optimal Experimental Design, vol. 184, Springer Science & Business Media, 2006.
- [7] L. Pronzato, Optimal experimental design and some related control problems, *Automatica* 44 (2008) 303–325.
- [8] A. Attia, E. Constantinescu, Optimal experimental design for inverse problems in the presence of observation correlations, *SIAM J. Sci. Comput.* 44 (2022) A2808–A2842.
- [9] N. Aretz-Nellesen, P. Chen, M.A. Grepl, K. Veroy, A sequential sensor selection strategy for hyper-parameterized linear Bayesian inverse problems, in: *Numerical Mathematics and Advanced Applications, ENUMATH 2019*, Springer, 2021, pp. 489–497.
- [10] J. Hart, M. Gulian, I. Manickam, L.P. Swiler, Solving high-dimensional inverse problems with auxiliary uncertainty via operator learning with limited data, *J. Mach. Learn. Model. Comput.* (2023).
- [11] V. Kolehmainen, T. Tarvainen, S.R. Arridge, J.P. Kaipio, Marginalization of uninteresting distributed parameters in inverse problems-application to diffuse optical tomography, *Int. J. Uncertain. Quantificat.* 1 (2011).
- [12] R. Nicholson, N. Petra, J.P. Kaipio, Estimation of the Robin coefficient field in a Poisson problem with uncertain conductivity field, *Inverse Probl.* 34 (2018) 115005.
- [13] A. Alexanderian, N. Petra, G. Stadler, I. Sunseri, Optimal design of large-scale Bayesian linear inverse problems under reducible model uncertainty: good to know what you don't know, *SIAM/ASA J. Uncertain. Quantificat.* 9 (2021) 163–184.
- [14] A. Bartuska, L. Espath, R. Tempone, Small-noise approximation for Bayesian optimal experimental design with nuisance uncertainty, *Comput. Methods Appl. Mech. Eng.* 399 (2022) 115320.
- [15] C. Feng, Y.M. Marzouk, A layered multiple importance sampling scheme for focused optimal Bayesian experimental design, *arXiv preprint*, arXiv:1903.11187, 2019.
- [16] D. Uciński, E-optimum sensor selection for estimation of subsets of parameters, *Measurement* 187 (2022) 110286.
- [17] K. Koval, A. Alexanderian, G. Stadler, Optimal experimental design under irreducible uncertainty for linear inverse problems governed by PDEs, *Inverse Probl.* 36 (2020) 75007.
- [18] N. Aretz-Nellesen, M.A.M. Grepl, K. Veroy, 3D-VAR for parameterized partial differential equations: a certified reduced basis approach, *Adv. Comput. Math.* 45 (2019) 2369–2400.
- [19] P. Binev, A. Cohen, O. Mula, J. Nichols, Greedy algorithms for optimal measurements selection in state estimation using reduced models, *SIAM/ASA J. Uncertain. Quantificat.* 6 (2018) 1101–1126.
- [20] Y. Maday, A.T. Patera, J.D. Penn, M. Yano, A parameterized-background data-weak approach to variational data assimilation: formulation, analysis, and application to acoustics, *Int. J. Numer. Methods Eng.* 102 (2015) 933–965.
- [21] M. Barrault, Y. Maday, N.C. Nguyen, A.T. Patera, An 'empirical interpolation' method: application to efficient reduced-basis discretization of partial differential equations, *C. R. Math.* 339 (2004) 667–672.
- [22] Y. Maday, O. Mula, A generalized empirical interpolation method: application of reduced basis techniques to data assimilation, in: *Analysis and Numerics of Partial Differential Equations*, Springer, 2013, pp. 221–235.
- [23] A. Alexanderian, P.J. Gloor, O. Ghattas, On Bayesian A- and D-optimal experimental designs in infinite dimensions, *Bayesian Anal.* 11 (2016) 671–695.
- [24] A. Alexanderian, N. Petra, G. Stadler, O. Ghattas, A-optimal design of experiments for infinite-dimensional Bayesian linear inverse problems with regularized ℓ_0 -sparsification, *SIAM J. Sci. Comput.* 36 (2014) A2122–A2148.
- [25] A. Attia, A. Alexanderian, A.K. Saibaba, Goal-oriented optimal design of experiments for large-scale Bayesian linear inverse problems, *Inverse Probl.* 34 (2018) 095009.
- [26] A. Alexanderian, A.K. Saibaba, Efficient D-optimal design of experiments for infinite-dimensional Bayesian linear inverse problems, *SIAM J. Sci. Comput.* 40 (2018) A2956–A2985.
- [27] K. Wu, P. Chen, O. Ghattas, An offline-online decomposition method for efficient linear Bayesian goal-oriented optimal experimental design: Application to optimal sensor placement, *SIAM J. Sci. Comput.* 45 (2023) B57–B77.

- [28] A. Alexanderian, Optimal experimental design for infinite-dimensional Bayesian inverse problems governed by PDEs: a review, *Inverse Probl.* (2021).
- [29] A. Alexanderian, N. Petra, G. Stadler, O. Ghattas, A fast and scalable method for A-optimal design of experiments for infinite-dimensional Bayesian nonlinear inverse problems, *SIAM J. Sci. Comput.* 38 (2016) A243–A272.
- [30] X. Huan, Y.M. Marzouk, Simulation-based optimal Bayesian experimental design for nonlinear systems, *J. Comput. Phys.* 232 (2013) 288–317.
- [31] K. Wu, P. Chen, O. Ghattas, A fast and scalable computational framework for large-scale and high-dimensional Bayesian optimal experimental design, *SIAM/AMS J. Uncertain. Quantificat.* 11 (2023) 235–261.
- [32] K. Wu, T. O’Leary-Roseberry, P. Chen, O. Ghattas, Large-scale Bayesian optimal experimental design with derivative-informed projected neural network, *J. Sci. Comput.* 95 (2023) 30.
- [33] T. Bui-Thanh, O. Ghattas, J. Martin, G. Stadler, A computational framework for infinite-dimensional Bayesian inverse problems part I: the linearized case, with application to global seismic inversion, *SIAM J. Sci. Comput.* 35 (2013) A2494–A2523.
- [34] T. Cui, Y. Marzouk, K. Willcox, Scalable posterior approximations for large-scale Bayesian inverse problems via likelihood-informed parameter and state reduction, *J. Comput. Phys.* 315 (2016) 363–387.
- [35] M.T. Parente, J. Wallin, B. Wohlmuth, Generalized bounds for active subspaces, *Electron. J. Stat.* 14 (2020) 917–943.
- [36] C. Lieberman, K. Willcox, O. Ghattas, Parameter and state model reduction for large-scale statistical inverse problems, *SIAM J. Sci. Comput.* 32 (2010) 2523–2542.
- [37] P. Chen, O. Ghattas, Hessian-based sampling for high-dimensional model reduction, *Int. J. Uncertain. Quantificat.* 9 (2019).
- [38] P. Chen, K. Wu, J. Chen, T. O’Leary-Roseberry, O. Ghattas, Projected Stein variational Newton: a fast and scalable Bayesian inference method in high dimensions, *NeurIPS* (2019), <https://arxiv.org/abs/1901.08659>.
- [39] P. Chen, O. Ghattas, Projected Stein variational gradient descent, in: *Advances in Neural Information Processing Systems*, 2020.
- [40] O. Zahm, T. Cui, K. Law, A. Spantini, Y. Marzouk, Certified dimension reduction in nonlinear Bayesian inverse problems, *Math. Comput.* 91 (2022) 1789–1835.
- [41] E. Qian, M. Grepl, K. Veroy, K. Willcox, A certified trust region reduced basis approach to PDE-constrained optimization, *SIAM J. Sci. Comput.* 39 (2017) S434–S460.
- [42] P. Chen, Model order reduction techniques for uncertainty quantification problems, Technical Report, 2014.
- [43] P. Chen, A. Quarteroni, G. Rozza, Reduced basis methods for uncertainty quantification, *SIAM/ASA J. Uncertain. Quantificat.* 5 (2017) 813–869.
- [44] T. O’Leary-Roseberry, U. Villa, P. Chen, O. Ghattas, Derivative-informed projected neural networks for high-dimensional parametric maps governed by PDEs, *Comput. Methods Appl. Mech. Eng.* 388 (2022) 114199.
- [45] T. O’Leary-Roseberry, P. Chen, U. Villa, O. Ghattas, Derivate informed neural operator: an efficient framework for high-dimensional parametric derivative learning, *arXiv:2206.10745*, 2022.
- [46] J.B. Nagel, Bayesian techniques for inverse uncertainty quantification, *IBK Ber.* 504 (2019).
- [47] G. Da Prato, *An Introduction to Infinite-Dimensional Analysis*, Springer Science & Business Media, 2006.
- [48] T. Cui, Y.M. Marzouk, K.E. Willcox, Data-driven model reduction for the Bayesian solution of inverse problems, *Int. J. Numer. Methods Eng.* 102 (2015) 966–990.
- [49] T. Bui-Thanh, K. Willcox, O. Ghattas, Model reduction for large-scale systems with high-dimensional parametric input space, *SIAM J. Sci. Comput.* 30 (2008) 3270–3288.
- [50] C. Schwab, R. Stevenson, Space-time adaptive wavelet methods for parabolic evolution problems, *Math. Comput.* 78 (2009) 1293–1318.
- [51] Q. Long, M. Scavino, R. Tempone, S. Wang, Fast estimation of expected information gains for Bayesian experimental designs based on Laplace approximations, *Comput. Methods Appl. Mech. Eng.* 259 (2013) 24–39.
- [52] N. Aretz, Data Assimilation and Sensor selection for Configurable Forward Models: Challenges and Opportunities for Model Order Reduction Methods, Ph.D. thesis, RWTH Aachen University, 2022.
- [53] D.E. Kirk, *Optimal Control Theory: an Introduction*, Courier Corporation, 2004.
- [54] M. Yano, J.D. Penn, A.T. Patera, A model-data weak formulation for simultaneous estimation of state and model bias, *C. R. Math.* 351 (2013) 937–941.
- [55] A. Cohen, W. Dahmen, R. DeVore, J. Fadili, O. Mula, J. Nichols, Optimal reduced model algorithms for data-based state estimation, *SIAM J. Numer. Anal.* 58 (2020) 3355–3381.
- [56] P. Benner, S. Gugercin, K. Willcox, A survey of projection-based model reduction methods for parametric dynamical systems, *SIAM Rev.* 57 (2015) 483–531.
- [57] W.H.A. Schilders, H.A. Van der Vorst, J. Rommes, *Model Order Reduction: Theory, Research Aspects and Applications*, vol. 13, Springer, 2008.
- [58] J.S. Hesthaven, G. Rozza, B. Stamm, *Certified Reduced Basis Methods for Parametrized Partial Differential Equations*, vol. 590, Springer, 2016.
- [59] A. Quarteroni, A. Manzoni, F. Negri, *Reduced Basis Methods for Partial Differential Equations: an Introduction*, vol. 92, Springer, 2015.
- [60] B. Haasdonk, Reduced basis methods for parametrized PDEs—a tutorial introduction for stationary and instationary problems, in: *Model Reduction and Approximation: Theory and Algorithms*, vol. 15, 2017, p. 65.
- [61] P. Binev, A. Cohen, W. Dahmen, R. DeVore, G. Petrova, P. Wojtaszczyk, Convergence rates for greedy algorithms in reduced basis methods, *SIAM J. Math. Anal.* 43 (2011) 1457–1472.
- [62] A. Buffa, Y. Maday, A.T. Patera, C. Prud’homme, G. Turinici, A priori convergence of the greedy algorithm for the parametrized reduced basis method, *ESAIM: Math. Model. Numer. Anal. (Modélisation Mathématique et Analyse Numérique)* 46 (2012) 595–603.
- [63] J. Jagalur-Mohan, Y.M. Marzouk, Batch greedy maximization of non-submodular functions: guarantees and applications to experimental design, *J. Mach. Learn. Res.* 22 (2021) 251–252.
- [64] J.L. Eftang, A.T. Patera, E.M. Rønquist, An “hp” certified reduced basis method for parametrized elliptic partial differential equations, *SIAM J. Sci. Comput.* 32 (2010) 3170–3200.
- [65] J.L. Eftang, D.J. Knezevic, A.T. Patera, An hp certified reduced basis method for parametrized parabolic partial differential equations, *Math. Comput. Model. Dyn. Syst.* 17 (2011) 395–422.
- [66] G.H. Golub, C.F. Van Loan, *Matrix Computations*, vol. 3, JHU Press, 2013.
- [67] K. Regenauer-Lieb, F. Horowitz, The Perth Basin geothermal opportunity, *Pet. West. Aust.* 3 (2007).
- [68] S. Corbel, O. Schilling, F.G. Horowitz, L.B. Reid, H.A. Sheldon, N.E. Timms, P. Wilkes, Identification and geothermal influence of faults in the Perth metropolitan area, Australia, in: *Thirty-Seventh Workshop on Geothermal Reservoir Engineering*, Stanford, CA, 2012.
- [69] H.A. Sheldon, B. Florio, M.G. Trefry, L.B. Reid, L.P. Ricard, K.A.R. Ghorri, The potential for convection and implications for geothermal energy in the Perth Basin, Western Australia, *Hydrogeol. J.* 20 (2012) 1251–1268.
- [70] O. Schilling, H.A. Sheldon, L.B. Reid, S. Corbel, Hydrothermal models of the Perth metropolitan area, Western Australia: implications for geothermal energy, *Hydrogeol. J.* 21 (2013) 605–621.
- [71] M. Pujol, L.P. Ricard, G. Bolton, 20 years of exploitation of the Yarragadee aquifer in the Perth Basin of Western Australia for direct-use of geothermal heat, *Geothermics* 57 (2015) 39–55.
- [72] J.F. Wellmann, L.B. Reid, Basin-scale geothermal model calibration: experience from the Perth Basin, Australia, *Energy Proc.* 59 (2014) 382–389.
- [73] D. Degen, K. Veroy, F. Wellmann, Certified reduced basis method in geosciences, *Comput. Geosci.* 24 (2020) 241–259.
- [74] D.M. Degen, Application of the reduced basis method in geophysical simulations: concepts, implementation, advantages, and limitations, Dissertation, RWTH Aachen University, 2020.
- [75] M. Bauer, W. Freeden, H. Jacobi, T. Neu, *Handbuch Tiefe Geothermie*, Springer, 2014.
- [76] J. Freymark, J. Sippel, M. Scheck-Wenderoth, K. Bär, M. Stiller, J.-G. Fritsche, M. Kracht, The deep thermal field of the upper Rhine graben, *Tectonophysics* 694 (2017) 114–129.

- [77] D. Degen, K. Veroy, J. Freymark, M. Scheck-Wenderoth, T. Poulet, F. Wellmann, Global sensitivity analysis to optimize basin-scale conductive model calibration—a case study from the upper Rhine graben, *Geothermics* 95 (2021) 102143.
- [78] M. de la Varga, A. Schaaf, F. Wellmann, GemPy 1.0: open-source stochastic geological modeling and inversion, *Geosci. Model Dev.* 12 (2019) 1–32.
- [79] C.J. Permann, D.R. Gaston, D. Andrš, R.W. Carlsen, F. Kong, A.D. Lindsay, J.M. Miller, J.W. Peterson, A.E. Slaughter, R.H. Stogner, MOOSE: enabling massively parallel multiphysics simulation, *SoftwareX* 11 (2020) 100430.
- [80] W. Dahmen, C. Plesken, G. Welper, Double greedy algorithms: reduced basis methods for transport dominated problems, *ESAIM: Math. Model. Numer. Anal. (Modélisation Mathématique et Analyse Numérique)* 48 (2014) 623–663.
- [81] D. Deming, Application of bottom-hole temperature corrections in geothermal studies, *Geothermics* 18 (1989) 775–786.
- [82] C. Vivas, S. Salehi, J.D. Tuttle, B. Rickard, Challenges and opportunities of geothermal drilling for renewable energy generation, *GRC Trans.* 44 (2020) 904–918.
- [83] N. Cressie, The origins of Kriging, *Math. Geol.* 22 (1990) 239–252.
- [84] F.L. Holgate, E.J. Gerner, OzTemp Well temperature data, Geoscience Australia, 2010, <http://www.ga.gov.au>, Catalogue.
- [85] Y. Maday, T. Anthony, J.D. Penn, M. Yano, PBDW state estimation: noisy observations; configuration-adaptive background spaces; physical interpretations, *ESAIM Proc. Surv.* 50 (2015) 144–168.
- [86] T. Taddei, Model order reduction methods for data assimilation: state estimation and structural health monitoring, Ph.D. thesis, Massachusetts Institute of Technology, 2017.