

# Metrics for Bayesian Optimal Experiment Design Under Model Misspecification

Tommie A. Catanach & Niladri Das

**Abstract**—The conventional approach to Bayesian decision-theoretic experiment design involves searching over possible experiments to select a design that maximizes the expected value of a specified utility function. The expectation is over the joint distribution of all unknown variables implied by the statistical model that will be used to analyze the collected data. Utility functions define experiments’ objectives; a common utility function is information gain. This article introduces an expanded framework for experimental design, where we go beyond the traditional Expected Information Gain criteria. We introduce Expected General Information Gain which measures robustness to the model discrepancy, and Expected Discriminatory Information to quantify how well an experiment can detect model discrepancy. The functionality of the framework is showcased through its application to a scenario involving a linearized spring mass damper system and an F-16 model where the model discrepancy is taken into account while doing Bayesian optimal experiment design.

## I. INTRODUCTION

For science and engineering systems there are often many choices of experiments to run, or data to collect, in order to infer information. Each of these choices has different costs in terms of time, money, or other resources. A common approach to designing the experiment stems from the field of Bayesian optimal experimental design (BOED). This approach uses the rigor of the Bayesian paradigm and information theory to formalize the design of experiments and treats it as an optimization problem. The aim is to maximize a utility function that captures the value of a particular experimental design. This utility function, typically the Expected Information Gain (EIG), depends on the posterior distribution sampled over many hypothetical realizations of plausible datasets from the experiment. However, for real applications, where there is model discrepancy, EIG is not the only measure of information we should consider.

In this work, we introduce two additional criteria that measure notions of robustness of a design. The first criterion, Expected Generalized Information Gain (EGIG), captures the expected information gained (or lost) when an experimenter uses a model with discrepancy. The second criterion, Expected Discriminatory Information (EDI) reflects whether the information gained from an experiment would be sufficient to discriminate between the model and an alternative. The EGIG-based design seeks to mitigate discrepancy while the

EDI-based one seeks to only detect it. With these criteria, we aim to correct pathological issues in BOED and advance the BOED literature, which has only a few works concerning the robustness of BOED.

In [1] a Bayesian linear regression example is shown where the system is analysed without considering model discrepancies. There, not only is the parameter under-estimated but the posterior credible intervals are not even close to covering the true parameter value, which is alarming. In practice, despite the theoretical elegance and optimal performance for accurate models, BOED may encounter significant issues if our model is not properly specified. This means that there is no value of  $x^*$  for which  $p(y|x = x^*, d)$  corresponds to the true distribution for  $p(y|d)$ , as noted in references [2] and [3]. Although model misspecification is a common problem in Bayesian settings, BOED methods are especially vulnerable because they use the model not only to fit data, but also to generate new data. The main issue is that Bayesian approaches only account for uncertainty in the model parameters, not in the model’s correctness, which can lead to disastrous outcomes where BOED continuously queries similar designs and produces low-quality datasets. Eliminating misspecification entirely is unrealistic, particularly in a general BOED context. However, a lot of work can be done to improve our comprehension and management of it. Presently, there is only a limited amount of research that covers both the theoretical [4],[5],[6], and [7] and empirical implications of misspecification [8], and very little has been done to examine the specific mechanisms that can lead to failures. This is where our EGIG and EDI metrics play an important role in evaluating the model robustness and identifying modeling failures. Some Bayesian-adjacent approaches that call out the need for robustness and optimality in design are [9] and [10]. Most notably, [9] considers robust sensor placement for linear dynamical systems under asymptotic D-optimal design.

*Outline:* Section II introduces the model and key concepts, Section III presents the BOED criteria, Section IV studies EGIG and EDI for two examples systems, and Section V provides discussions.

## II. MODELING AND KEY CONCEPTS

### A. System Description

We will study BOED in the context of simplified models, specifically stationary discrete-time linear processes driven

Tommie A. Catanach (tacatan@sandia.gov) and Niladri Das (corresponding author, ndas@sandia.gov), are with Sandia National Laboratories, Livermore, CA 94550, USA.

by Gaussian noise. We define the state vector as  $\mathbf{x}_t \in \mathbb{R}^n$ ,

$$\mathbf{x}_t = \mathbf{A}\mathbf{x}_{t-1} + \boldsymbol{\eta}_t, \quad t = 1, 2, \dots \quad (1)$$

$\mathbf{A}$  is an  $n \times n$  transition matrix and  $\boldsymbol{\eta}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q})$  is the process noise where  $\mathbf{Q} \succeq \mathbf{0}$ . We assume  $\mathbf{x}_0 \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ . For simplicity, unless specified, we will take  $\boldsymbol{\mu}_0 = \mathbf{0}$ .

The observation equation is,

$$\mathbf{y}_t = \mathbf{H}\mathbf{x}_t + \mathbf{v}_t, \quad (2)$$

where the measurements are  $\mathbf{y}_t \in \mathbb{R}^s$ ,  $\mathbf{H}$  is the measurement matrix and  $\mathbf{v}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$  where  $\mathbf{R} \succ \mathbf{0}$ . The random vectors  $\{\mathbf{x}_0, \boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_t, \mathbf{v}_1, \dots, \mathbf{v}_t\}$  are all assumed to be independent.

From this general case, we will study two simplifications. First, we consider a system without dynamics (or equivalently a single time step of the system), corresponding to  $\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{v}$ . Second, we will study the system after it has converged to its stationary distribution, assuming that  $\mathbf{A}$  is asymptotically stable. In this case, if  $t$  is sufficiently large, we have that  $\mathbf{x}_t \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_L)$ , where  $\boldsymbol{\Sigma}_L$  is the solution to the discrete Lyapunov equation,  $\boldsymbol{\Sigma}_L = \mathbf{A}\boldsymbol{\Sigma}_L\mathbf{A}^T + \mathbf{Q}$ .

### B. Bayesian Inference

In Bayesian inference, to rigorously update our beliefs about  $\mathbf{X}$  with observation data  $\mathbf{Y}$ , we apply Bayes' theorem,

$$p(\mathbf{X} | \mathbf{Y}) = \frac{p(\mathbf{Y} | \mathbf{X})p(\mathbf{X})}{p(\mathbf{Y})}. \quad (3)$$

The prior  $p(\mathbf{X})$  reflects our initial beliefs about  $\mathbf{X}$  while  $p(\mathbf{X} | \mathbf{Y})$  is our posterior (after observations) belief. The likelihood,  $p(\mathbf{Y} | \mathbf{X})$  is the probability of observing  $\mathbf{Y}$  given a state  $\mathbf{X}$ , while  $p(\mathbf{Y})$  is the overall probability of observing the data given our prior (called the evidence). Often we are interested in measuring how informative is the data. To do this we measure our change in belief, i.e., the information gain, using the Kullback–Leibler (KL) divergence,

$$D_{\text{KL}}[p(\mathbf{X} | \mathbf{Y}) || p(\mathbf{X})] = \int p(\mathbf{X} | \mathbf{Y}) \log \frac{p(\mathbf{X} | \mathbf{Y})}{p(\mathbf{X})} d\mathbf{X} \quad (4)$$

For the Gaussian case where  $p(\mathbf{X} | \mathbf{Y}) \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$  and  $p(\mathbf{X}) \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$  the KL divergence is,

$$\frac{1}{2} \left( \text{Tr}[\boldsymbol{\Sigma}_0^{-1}\boldsymbol{\Sigma}_1] - n + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}_0^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) + \log \frac{|\boldsymbol{\Sigma}_0|}{|\boldsymbol{\Sigma}_1|} \right). \quad (5)$$

The KL divergence can be generalized using a more expressive, yet still information theoretically valid, measure of information [11] defined over three distributions,  $r(\mathbf{X})$ ,  $p(\mathbf{X})$ , and  $q(\mathbf{X})$ , given by:

$$\mathcal{I}_{r(\mathbf{X})}[p(\mathbf{X}) || q(\mathbf{X})] = \int r(\mathbf{X}) \log \frac{p(\mathbf{X})}{q(\mathbf{X})} d\mathbf{X}. \quad (6)$$

The interpretation of this form of information is that we want to measure a change in belief (e.g., information gained or lost) when updating from  $q(\mathbf{X})$  to  $p(\mathbf{X})$  in the view of  $r(\mathbf{X})$ . The view defines our reference frame for assessing changes in information. Typically, both  $r(\mathbf{X})$  and  $p(\mathbf{X})$  would represent the posterior with  $q(\mathbf{X})$  as the prior, thus recovering the regular KL divergence expression. However, in the case where there is model discrepancy,  $r(\mathbf{X})$  could be the unknown posterior from the true model, while  $p(\mathbf{X})$  could be the inferred posterior from the model with discrepancy. Therefore, we could measure whether inference with the model discrepancy is still getting close to the correct result. We note that unlike KL divergence this measure can be negative, meaning that  $q(\mathbf{X})$  provides more information about  $r(\mathbf{X})$  than  $p(\mathbf{X})$  does.

For the case where,  $r(\mathbf{X})$ ,  $p(\mathbf{X})$ , and  $q(\mathbf{X})$  are all described by multivariate Gaussians,

$$\begin{aligned} \mathcal{I}_{r(\mathbf{X})}[p(\mathbf{X}) || q(\mathbf{X})] = & \frac{1}{2} \left( \text{Tr}[(\boldsymbol{\Sigma}_q^{-1} - \boldsymbol{\Sigma}_p^{-1})\boldsymbol{\Sigma}_r] - (\boldsymbol{\mu}_r - \boldsymbol{\mu}_p)^T \boldsymbol{\Sigma}_p^{-1} (\boldsymbol{\mu}_r - \boldsymbol{\mu}_p) \right. \\ & \left. + (\boldsymbol{\mu}_r - \boldsymbol{\mu}_q)^T \boldsymbol{\Sigma}_q^{-1} (\boldsymbol{\mu}_r - \boldsymbol{\mu}_q) + \log \frac{|\boldsymbol{\Sigma}_q|}{|\boldsymbol{\Sigma}_p|} \right). \quad (7) \end{aligned}$$

This uses the fact that eq. 6 can be expressed as the difference of two KL divergences and employ eq. 5.

### C. Bayesian Filtering

For a Markov process where the state  $\mathbf{x}_t$  only depends on  $\mathbf{x}_{t-1}$  and the observation  $\mathbf{y}_t$  only depends on  $\mathbf{x}_t$  we can simplify the inference problem for the state  $\mathbf{x}_t$  given a time series of observations  $\mathbf{Y}_t = \{\mathbf{y}_0 \dots \mathbf{y}_t\}$  as

$$p(\mathbf{x}_t | \mathbf{Y}_t) = \frac{p(\mathbf{y}_t | \mathbf{x}_t)p(\mathbf{x}_t | \mathbf{Y}_{t-1})}{p(\mathbf{Y} | \mathbf{Y}_{t-1})}. \quad (8)$$

Using this, the Bayesian filter for the system described by eq.1-2, is the Kalman filter,

$$\boldsymbol{\mu}_{t|t-1} = \mathbf{A}\boldsymbol{\mu}_{t-1|t-1} \quad (9)$$

$$\boldsymbol{\Sigma}_{t|t-1} = \mathbf{A}\boldsymbol{\Sigma}_{t-1|t-1}\mathbf{A}^T + \mathbf{Q} \quad (10)$$

$$\boldsymbol{\mu}_{t|t} = \boldsymbol{\mu}_{t|t-1} + \mathbf{K}_t(\mathbf{y}_t - \mathbf{H}\boldsymbol{\mu}_{t|t-1}) \quad (11)$$

$$\boldsymbol{\Sigma}_{t|t} = (\mathbf{I} - \mathbf{K}_t\mathbf{H})\boldsymbol{\Sigma}_{t|t-1}, \quad (12)$$

where  $\mathbf{K}_t = \boldsymbol{\Sigma}_{t|t-1}\mathbf{H}^T\mathbf{S}_t^{-1}$  is the Kalman gain matrix,  $\mathbf{S}_t = \mathbf{H}\boldsymbol{\Sigma}_{t|t-1}\mathbf{H}^T + \mathbf{R}$  is the predictive uncertainty, and  $\mathbf{I}$  is the identity matrix. Considering a single time step, the *a-priori* estimator of  $\mathbf{x}_t$  is  $\boldsymbol{\mu}_{t|t-1}$  with covariance  $\boldsymbol{\Sigma}_{t|t-1}$ . The *a-posteriori* estimator of  $\mathbf{x}_t$  is  $\boldsymbol{\mu}_{t|t}$  with covariance  $\boldsymbol{\Sigma}_{t|t}$ . Therefore, the prior, posterior, and evidence are:

$$p(\mathbf{x}_t) \sim \mathcal{N}(\boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}), \quad (13)$$

$$p(\mathbf{x}_t | \mathbf{y}_t, \mathbf{d}) \sim \mathcal{N}(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}), \quad (14)$$

$$p(\mathbf{y}_t | \mathbf{d}) \sim \mathcal{N}(\mathbf{H}\boldsymbol{\mu}_{t|t-1}, \mathbf{S}_t). \quad (15)$$

As we can see from eq. 9-12, only the means  $\boldsymbol{\mu}$  depend on the observations  $\mathbf{y}$ . Thus, when  $\mathbf{A}$  is asymptotically stable we can find the stationary distribution of  $\boldsymbol{\Sigma}_{t|t}$ . We define  $\boldsymbol{\Sigma}_{t|t-1} \rightarrow \boldsymbol{\Gamma}$  and  $\boldsymbol{\Sigma}_{t|t} \rightarrow \boldsymbol{\Sigma}_D$  as  $t \rightarrow \infty$ . We first use the discrete time algebraic Riccati equation (DARE) to give,

$$\boldsymbol{\Gamma} = \mathbf{A}\boldsymbol{\Gamma}\mathbf{A}^T + \mathbf{Q} - \mathbf{A}\boldsymbol{\Gamma}\mathbf{H}^T(\mathbf{H}\boldsymbol{\Gamma}\mathbf{H}^T + \mathbf{R})^{-1}\mathbf{H}\boldsymbol{\Gamma}\mathbf{A}^T, \quad (16)$$

and then solve for  $\boldsymbol{\Sigma}_D$  via:

$$\boldsymbol{\Sigma}_D = \boldsymbol{\Gamma} - \boldsymbol{\Gamma}\mathbf{H}^T(\mathbf{H}\boldsymbol{\Gamma}\mathbf{H}^T + \mathbf{R})^{-1}\mathbf{H}\boldsymbol{\Gamma}. \quad (17)$$

#### D. Bayesian Optimal Experimental Design

In BOED, the first step to modeling the problem is to define a utility function  $U(\mathbf{d})$  that gives the value of performing an experiment at  $\mathbf{d} \in \mathcal{D}$ . The set  $\mathcal{D}$  spans the space of possible designs. In Bayesian design, the utility is a function of the posterior distribution  $p(\mathbf{X} | \mathbf{d}, \mathbf{Y})$ . The utility function is maximized to find the optimal design  $\mathbf{d}^*$ , i.e.  $\mathbf{d}^* = \operatorname{argmax}_{\mathbf{d} \in \mathcal{D}} U(\mathbf{d})$ . The choice of the utility function  $U(\mathbf{d})$  is crucial, as different functions will usually lead to different optimal designs [12]. A principled choice often used in BOED is mutual information. This is the information gained about  $\mathbf{X}$  by taking measurements,  $\mathbf{Y}$ , according to design  $\mathbf{d}$ . This is just the KL divergence from prior to posterior,  $\mathbb{D}_{\text{KL}}[p(\mathbf{X} | \mathbf{Y}, \mathbf{d}) || p(\mathbf{X})]$ , eq. 4.

However at the point of choosing  $\mathbf{d}$ , we do not have measurements. So to assess the effectiveness of design  $\mathbf{d}$ , we take the expected KL divergence over plausible datasets  $p(\mathbf{Y} | \mathbf{d})$ . This utility function is known as Expected Information Gain (EIG) and is defined as,

$$\begin{aligned} \text{EIG}(\mathbf{d}) &= \mathbb{E}_{p(\mathbf{Y} | \mathbf{d})} \left[ \mathbb{D}_{\text{KL}}(p(\mathbf{X} | \mathbf{Y}, \mathbf{d}) || p(\mathbf{X})) \right] \\ &= \int p(\mathbf{X}, \mathbf{Y} | \mathbf{d}) \log \frac{p(\mathbf{X} | \mathbf{Y}, \mathbf{d})}{p(\mathbf{X})} d\mathbf{X} d\mathbf{Y}. \end{aligned} \quad (18)$$

### III. DESIGN CRITERIA

#### A. Expected Information Gain

For the linear Gaussian model given by eq.1-2, we can derive expressions for the EIG.

*Single Step Update:* First, for the case of a single update step (or equivalently, no dynamics) we begin by substituting the values from eq. 9 - 12 into the Gaussian KL divergence expression, eq. 5. Rearranging terms with the matrix inversion lemma and cyclic property of the trace, the information gain from prior to posterior is

$$\begin{aligned} \mathbb{D}_{\text{KL}}[p(\mathbf{x}_t | \mathbf{y}_t, \mathbf{d}) || p(\mathbf{x}_t)] &= \\ \frac{1}{2} \left[ \log |\mathbf{I} + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H} \boldsymbol{\Sigma}_{t|t-1}| - \text{Tr}[\mathbf{S}_t^{-1} \mathbf{H} \boldsymbol{\Sigma}_{t|t-1} \mathbf{H}^T] \right. \\ &\quad \left. + (\mathbf{y}_t - \mathbf{H} \boldsymbol{\mu}_{t|t-1})^T \mathbf{S}_t^{-1} \mathbf{H} \boldsymbol{\Sigma}_{t|t-1} \mathbf{H}^T \mathbf{S}_t^{-1} (\mathbf{y}_t - \mathbf{H} \boldsymbol{\mu}_{t|t-1}) \right] \end{aligned} \quad (19)$$

Only the last term depends on  $\mathbf{y}_t$ , so for EIG we just need

to find the expectation of the quadratic term, which is,

$$\begin{aligned} &\mathbb{E}_{p(\mathbf{y}_t | \mathbf{d})} \left[ (\mathbf{y}_t - \mathbf{H} \boldsymbol{\mu}_{t|t-1})^T \mathbf{S}_t^{-1} \mathbf{H} \boldsymbol{\Sigma}_{t|t-1} \mathbf{H}^T \mathbf{S}_t^{-1} \right. \\ &\quad \left. (\mathbf{y}_t - \mathbf{H} \boldsymbol{\mu}_{t|t-1}) \right] \\ &= \text{Tr} \left[ \mathbf{S}_t^{-1} \mathbf{H} \boldsymbol{\Sigma}_{t|t-1} \mathbf{H}^T \mathbf{S}_t^{-1} \text{Cov}(\mathbf{y}_t - \mathbf{H} \boldsymbol{\mu}_{t|t-1}) \right] \\ &= \text{Tr} \left[ \mathbf{S}_t^{-1} \mathbf{H} \boldsymbol{\Sigma}_{t|t-1} \mathbf{H}^T \right]. \end{aligned} \quad (20)$$

Here we recall eq. 15 so  $(\mathbf{y}_t - \mathbf{H} \boldsymbol{\mu}_{t|t-1})$  has mean  $\mathbf{0}$  and covariance  $\mathbf{S}_t$ . Therefore, noting cancellation of trace terms, EIG of the single step of the Kalman filter is

$$\begin{aligned} \text{EIG}(\mathbf{d}) &= \mathbb{E}_{p(\mathbf{y}_t | \mathbf{d})} \left[ \mathbb{D}_{\text{KL}}(p(\mathbf{x}_t | \mathbf{y}_t, \mathbf{d}) || p(\mathbf{x}_t)) \right] \\ &= \frac{1}{2} \left[ \log |\mathbf{I} + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H} \boldsymbol{\Sigma}_{t|t-1}| \right]. \end{aligned} \quad (21)$$

*Infinite Horizon:* We may also be interested in assessing the EIG about a state  $\mathbf{x}_t$  when the system and filters have converged to their stationary distributions. For this, we define our prior knowledge about  $\mathbf{x}_t$  as the solution to the Lyapunov equation, e.g.,  $p(\mathbf{x}_t | \mathbf{d}) = \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_L)$ , when  $t$  is sufficiently large to be in the asymptotic regime. Similarly, when we have a sufficiently large set of observations,  $\mathbf{Y}_t$ , we know the posterior belief about  $\mathbf{x}_t$  will have the form  $p(\mathbf{x}_t | \mathbf{Y}_t, \mathbf{d}) = \mathcal{N}(\boldsymbol{\mu}_t(\mathbf{Y}_t), \boldsymbol{\Sigma}_D)$ . Here we express  $\boldsymbol{\mu}_t$  as a function of  $\mathbf{Y}_t$  to emphasize that  $\boldsymbol{\mu}_t$  is a random variable defined by  $\mathbf{Y}_t$ . Therefore, information gain from observing  $\mathbf{Y}_t$  is

$$\begin{aligned} \mathbb{D}_{\text{KL}}[p(\mathbf{x}_t | \mathbf{Y}_t, \mathbf{d}) || p(\mathbf{x}_t)] &= \\ \frac{1}{2} \left( \text{Tr} [\boldsymbol{\Sigma}_L^{-1} \boldsymbol{\Sigma}_D] - n + \boldsymbol{\mu}_t(\mathbf{Y}_t)^T \boldsymbol{\Sigma}_L^{-1} \boldsymbol{\mu}_t(\mathbf{Y}_t) + \log \frac{|\boldsymbol{\Sigma}_L|}{|\boldsymbol{\Sigma}_D|} \right). \end{aligned} \quad (22)$$

Again, only the quadratic term depends on observations. Therefore, to compute EIG we first derive the expectation,

$$\begin{aligned} &\mathbb{E}_{p(\mathbf{Y}_t | \mathbf{d})} \left[ \boldsymbol{\mu}_t(\mathbf{Y}_t)^T \boldsymbol{\Sigma}_L^{-1} \boldsymbol{\mu}_t(\mathbf{Y}_t) \right] \\ &= \text{Tr} \left[ \boldsymbol{\Sigma}_L^{-1} \text{Cov}(\boldsymbol{\mu}_t(\mathbf{Y}_t), \boldsymbol{\mu}_t(\mathbf{Y}_t)^T) \right] = \text{Tr} \left[ \boldsymbol{\Sigma}_L^{-1} (\boldsymbol{\Sigma}_L - \boldsymbol{\Sigma}_D) \right] \\ &= n - \text{Tr} \left[ \boldsymbol{\Sigma}_L^{-1} \boldsymbol{\Sigma}_D \right]. \end{aligned} \quad (23)$$

We use  $\mathbb{E}_{p(\mathbf{Y}_t | \mathbf{d})} [\boldsymbol{\mu}_t(\mathbf{Y}_t)] = \mathbf{0}$  and  $\mathbb{E} [\boldsymbol{\mu}_t(\mathbf{Y}_t) \boldsymbol{\mu}_t(\mathbf{Y}_t)^T] = (\boldsymbol{\Sigma}_L - \boldsymbol{\Sigma}_D)$ . This is shown as eq. 57 in Appendix VI-A.

Therefore, taking the expectation of eq. 22 over  $\mathbf{Y}_t$  and substituting in the result of eq. 23 which cancels the trace terms, we find similarly to eq. 21 that,

$$\begin{aligned} \text{EIG}(\mathbf{d}) &= \mathbb{E}_{p(\mathbf{Y}_t | \mathbf{d})} \left[ \mathbb{D}_{\text{KL}}(p(\mathbf{x}_t | \mathbf{Y}_t, \mathbf{d}) || p(\mathbf{x}_t)) \right] \\ &\rightarrow \frac{1}{2} \log \frac{|\boldsymbol{\Sigma}_L|}{|\boldsymbol{\Sigma}_D|}, \text{ as } t \rightarrow \infty. \end{aligned} \quad (24)$$

## B. Expected Generalized Information Gain

Using the generalized measure of information in eq. 6, we can assess how much information is expected to be gained or lost by an experiment  $d$  when there is model discrepancy. We define the true model as  $\mathcal{M}^*$  and the model with discrepancy as  $\mathcal{M}$ , both of which have the same unknown states,  $\mathbf{X}$ , which we seek to infer. This expectation is taken over data that is generated according to  $p(\mathbf{Y} | d, \mathcal{M}^*)$ . This leads to the Expected Generalized Information Gain (EGIG) given by,

$$\text{EGIG}(d, \mathcal{M}, \mathcal{M}^*) =$$

$$\mathbb{E}_{p(\mathbf{Y} | d, \mathcal{M}^*)} \left[ \mathcal{I}_{p^*} \left[ p(\mathbf{X} | \mathbf{Y}, d, \mathcal{M}) || p(\mathbf{X} | \mathcal{M}) \right] \right]$$

$$= \int p(\mathbf{X}, \mathbf{Y} | d, \mathcal{M}^*) \log \frac{p(\mathbf{X} | \mathbf{Y}, d, \mathcal{M})}{p(\mathbf{X} | \mathcal{M})} d\mathbf{X} d\mathbf{Y} \quad (25)$$

$$= \int p(\mathbf{X}, \mathbf{Y} | d, \mathcal{M}^*) \log \frac{p(\mathbf{Y} | \mathbf{X}, d, \mathcal{M})}{p(\mathbf{Y} | \mathcal{M})} d\mathbf{X} d\mathbf{Y} \quad (26)$$

For notation we have  $p^* := p(\mathbf{X} | \mathbf{Y}, d, \mathcal{M}^*)$ . Note that eq. 26 is a simple rearrangement using Bayes' theorem, which can be easier to compute for some problems.

In the following analysis we use  $\mathcal{M}^*$  as a theoretical quantity to derive our metric. Typically we do not know  $\mathcal{M}^*$ , so in practice we should either define a set of plausible models we want to be robust to or we can assess the sensitivity to perturbations away from  $\mathcal{M}$  by computing derivatives of the EGIG using either automatic differentiation or numerical derivatives. However, for some applications like surrogate modeling we do know  $\mathcal{M}^*$ .

In the context of inferring  $\mathbf{x}_t$  with a system defined by eq. 1-2 we define the true model  $\mathcal{M}^* = \{\mathbf{A}^*, \mathbf{H}^*, \mathbf{Q}^*, \mathbf{R}^*\}$  and the model we use for inference as  $\mathcal{M} = \{\mathbf{A}, \mathbf{H}, \mathbf{Q}, \mathbf{R}\}$ .

Single Step Update: We start with the EGIG of eq. 26. We defined  $\boldsymbol{\mu}_{t|t-1} = \mathbf{A}$ ,  $\boldsymbol{\mu}_{t|t-1}^* = \mathbf{A}^*$ ,  $\boldsymbol{\Sigma}_{t|t-1} = \mathbf{A}\boldsymbol{\Sigma}_{t-1|t-1}\mathbf{A}^T + \mathbf{Q}$ , and  $\boldsymbol{\Sigma}_{t|t-1}^* = \mathbf{A}^*\boldsymbol{\Sigma}_{t-1|t-1}^*\mathbf{A}^{*T} + \mathbf{Q}^*$ . We then note the distributions,

$$p(\mathbf{x}_t, \mathbf{y}_t | d, \mathcal{M}^*) =$$

$$\mathcal{N} \left( \begin{pmatrix} \boldsymbol{\mu}_{t|t-1}^* \\ \mathbf{H}^* \boldsymbol{\mu}_{t|t-1}^* \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{t|t-1}^* & \boldsymbol{\Sigma}_{t|t-1}^* \mathbf{H}^{*T} \\ \mathbf{H}^* \boldsymbol{\Sigma}_{t|t-1}^* & \mathbf{S}_t^* \end{pmatrix} \right) \quad (27)$$

$$p(\mathbf{y}_t | \mathbf{x}_t d, \mathcal{M}) = \mathcal{N}(\mathbf{H}\mathbf{x}_t, \mathbf{R}) \quad (28)$$

$$p(\mathbf{y}_t | d, \mathcal{M}) = \mathcal{N}(\mathbf{H}\boldsymbol{\mu}_{t|t-1}, \mathbf{S}_t). \quad (29)$$

Recall that  $\mathbf{S}_t = \mathbf{H}\boldsymbol{\Sigma}_{t|t-1}\mathbf{H}^T + \mathbf{R}$  and  $\mathbf{S}_t^* = \mathbf{H}^*\boldsymbol{\Sigma}_{t|t-1}^*\mathbf{H}^{*T} + \mathbf{R}^*$ . Substituting these distributions into eq. 26, we arrive at

$$\text{EGIG}(d, \mathcal{M}, \mathcal{M}^*) =$$

$$\mathbb{E}_{p(\mathbf{x}_t, \mathbf{y}_t | \mathcal{M}^*)} \left[ \log \frac{p(\mathbf{x}_t | \mathbf{y}_t, d, \mathcal{M})}{p(\mathbf{x}_t | d, \mathcal{M})} \right]$$

$$= \frac{1}{2} \left( \log \frac{|\mathbf{S}_t|}{|\mathbf{R}|} - \mathbb{E} \left[ (\mathbf{y}_t - \mathbf{H}\mathbf{x}_t)^T \mathbf{R}^{-1} (\mathbf{y}_t - \mathbf{H}\mathbf{x}_t) \right] \right)$$

$$+ \mathbb{E} \left[ (\mathbf{y}_t - \mathbf{H}\boldsymbol{\mu}_{t|t-1})^T \mathbf{S}_t^{-1} (\mathbf{y}_t - \mathbf{H}\boldsymbol{\mu}_{t|t-1}) \right]. \quad (30)$$

Using eq. 27-29, it is straight forward to compute the means and covariances of  $(\mathbf{y}_t - \mathbf{H}\mathbf{x}_t)$  and  $(\mathbf{y}_t - \mathbf{H}\boldsymbol{\mu}_{t|t-1})$ ,

$$(\mathbf{y}_t - \mathbf{H}\mathbf{x}_t) \sim \mathcal{N}((\mathbf{H}^* - \mathbf{H})\boldsymbol{\mu}_{t|t-1}^*,$$

$$(\mathbf{H}^* - \mathbf{H})^T \boldsymbol{\Sigma}_{t|t-1}^* (\mathbf{H}^* - \mathbf{H}) + \mathbf{R}^*) \quad (31)$$

$$(\mathbf{y}_t - \mathbf{H}\boldsymbol{\mu}_{t|t-1}) \sim \mathcal{N}(\mathbf{H}^* \boldsymbol{\mu}_{t|t-1}^* - \mathbf{H}\boldsymbol{\mu}_{t|t-1}, \mathbf{S}_t^*). \quad (32)$$

Given these distributions, it is useful to define the following variables  $\boldsymbol{\Delta}_H = \mathbf{H}^* - \mathbf{H}$  and  $\boldsymbol{\Delta}_y = (\mathbf{H}^* \boldsymbol{\mu}_{t|t-1}^* - \mathbf{H}\boldsymbol{\mu}_{t|t-1})$ . Therefore, we can define the EGIG as:

$$\text{EGIG}(d, \mathcal{M}, \mathcal{M}^*) =$$

$$\frac{1}{2} \left( \log \frac{|\mathbf{S}_t|}{|\mathbf{R}|} - \text{Tr}[\mathbf{R}^{-1} \boldsymbol{\Delta}_H \boldsymbol{\Sigma}_{t|t-1}^* \boldsymbol{\Delta}_H^T] - \text{Tr}[\mathbf{R}^{-1} \mathbf{R}^*] \right.$$

$$+ \text{Tr}[\mathbf{S}_t^{-1} \mathbf{S}_t^*] - \boldsymbol{\mu}_{t|t-1}^{*T} \boldsymbol{\Delta}_H^T \mathbf{R}^{-1} \boldsymbol{\Delta}_H \boldsymbol{\mu}_{t|t-1}^*$$

$$\left. + \boldsymbol{\Delta}_y^T \mathbf{S}_t^{-1} \boldsymbol{\Delta}_y \right). \quad (33)$$

Infinite Horizon: For the infinite horizon case for inferring  $\mathbf{x}_t$  we know that our prior and posterior are Gaussian. Therefore, when computing the EGIG we can use the expression in eq. 7 and then compute the expectation over observations  $\mathbf{Y}_t \sim p(\mathbf{Y}_t | \mathcal{M}^*)$ . Here,  $r(\mathbf{X})$  is  $p(\mathbf{x}_t | \mathbf{Y}_t, d, \mathcal{M}^*)$ ,  $p(\mathbf{X})$  is  $p(\mathbf{x}_t | \mathbf{Y}_t, d, \mathcal{M})$ , and  $q(\mathbf{X})$  is  $p(\mathbf{x}_t | \mathcal{M})$ . By inspection, we see again that the only terms that depend on  $\mathbf{Y}_t$  are the quadratic terms. Therefore, we begin with those terms.

First, we note the asymptotic results:  $\boldsymbol{\Sigma}_{t|t} \rightarrow \boldsymbol{\Sigma}_D$ ,  $\boldsymbol{\Sigma}_{t|t}^* \rightarrow \boldsymbol{\Sigma}_D^*$ ,  $\boldsymbol{\Sigma}_{t|0} \rightarrow \boldsymbol{\Sigma}_L$ ,  $\boldsymbol{\mu}_{t|0} = \mathbf{0}$ , and  $\boldsymbol{\mu}_{t|t}^* \stackrel{t \rightarrow \infty}{\sim} \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_L^* - \boldsymbol{\Sigma}_D^*)$ . This gives us:

$$\mathbb{E}_{p(\mathbf{Y}_t | \mathcal{M}^*)} \left[ \left( \boldsymbol{\mu}_{t|t}^* - \boldsymbol{\mu}_{t|0} \right)^T \boldsymbol{\Sigma}_L^{-1} \left( \boldsymbol{\mu}_{t|t}^* - \boldsymbol{\mu}_{t|0} \right) \right]$$

$$= \text{Tr} \left[ \boldsymbol{\Sigma}_L^{-1} (\boldsymbol{\Sigma}_L^* - \boldsymbol{\Sigma}_D^*) \right]. \quad (34)$$

For the second expectation, we again rely on results presented in detail in Appendix VI-A. First,  $\mathbb{E}_{p(\mathbf{Y}_t | \mathcal{M}^*)} [\boldsymbol{\mu}_{t|t}^* - \boldsymbol{\mu}_{t|t}] = \mathbf{0}$ . Second, therefore,

$$\mathbb{E}_{p(\mathbf{Y}_t | \mathcal{M}^*)} \left[ \left( \boldsymbol{\mu}_{t|t}^* - \boldsymbol{\mu}_{t|t} \right)^T \boldsymbol{\Sigma}_D^{-1} \left( \boldsymbol{\mu}_{t|t}^* - \boldsymbol{\mu}_{t|t} \right) \right]$$

$$= \text{Tr}[\boldsymbol{\Sigma}_D^{-1} \mathbf{M}_\Delta]. \quad (35)$$

$$\mathbf{M}_\Delta = \text{Cov}(\boldsymbol{\mu}_{t|t}^* - \boldsymbol{\mu}_{t|t}) = [-\mathbb{I} \ \mathbb{I}] \mathbf{M} [-\mathbb{I} \ \mathbb{I}]^T \quad (36)$$

where  $\mathbf{M}$ , the asymptotic covariance matrix of  $[\boldsymbol{\mu}_{t|t} \ \boldsymbol{\mu}_{t|t}^*]^T$ , is the solution to the Lyapunov equation given by:

$$\mathbf{M} = \mathbf{A}\mathbf{M}\mathbf{A}^T + \boldsymbol{\mathcal{Q}} \quad (37)$$

$$\mathcal{A} = \begin{pmatrix} (I - KH)A & KH^*A^* \\ \mathbf{0} & A^* \end{pmatrix} \quad (38)$$

$$\mathcal{Q} = \begin{pmatrix} KS^*K^T & KS^*K^{*T} \\ K^*S^*K^T & K^*S^*K^{*T} \end{pmatrix} \quad (39)$$

Therefore, using these two equations we arrive at the EGIG for the infinite horizon system,

$$\begin{aligned} \text{EGIG}(d, \mathcal{M}, \mathcal{M}^*) &= \mathbb{E}_{p(\mathbf{Y}_t | \mathcal{M}^*)} [ \\ &\mathcal{I}_{p(\mathbf{x}_t | \mathbf{Y}_t, d, \mathcal{M}^*)} [p(\mathbf{x}_t | \mathbf{Y}_t, d, \mathcal{M}) || p(\mathbf{x}_t | \mathcal{M})] \rightarrow \\ &\frac{1}{2} \left( \text{Tr}[\Sigma_L^{-1} \Sigma_L^*] - \text{Tr}[\Sigma_D^{-1} (\Sigma_D^* + M_\Delta)] + \log \frac{|\Sigma_L|}{|\Sigma_D|} \right) \\ &\text{as } t \rightarrow \infty. \end{aligned} \quad (40)$$

### C. Expected Discriminatory Information

While EIG measures efficiency and EGIG measures robustness, we introduce the Expected Discriminatory Information (EDI) criteria to quantify how well an experiment can identify modeling failures. As such, unlike EGIG which is focused on comparing the Bayesian inference solution in the domain of the states  $\mathbf{x}$ , EDI compares them in the data domain,  $\mathbf{y}$ . Therefore, we can compare models that have different states and forms, e.g., different number of states. The EDI takes inspiration from the use of Bayes factors to compare models. Therefore we define the EDI as the expected Bayes factor given data from a true model  $\mathcal{M}^*$ :

$$\begin{aligned} \text{EDI}(d, \mathcal{M}, \mathcal{M}^*) &= D_{\text{KL}} [p(\mathbf{Y} | d, \mathcal{M}^*) || p(\mathbf{Y} | d, \mathcal{M})] \\ &= \int p(\mathbf{Y} | d, \mathcal{M}^*) \log \frac{p(\mathbf{Y} | d, \mathcal{M}^*)}{p(\mathbf{Y} | d, \mathcal{M})} d\mathbf{Y}. \end{aligned} \quad (41)$$

For filtering where  $\mathbf{Y}_t = \{y_0 \dots y_t\}$ , we can express the EDI using an iterative update leveraging a similar strategy for computing model evidence using a Bayesian filter,

$$\begin{aligned} \text{EDI}(d, \mathcal{M}, \mathcal{M}^*, t) &= \mathbb{E}_{p(\mathbf{y}_t, \mathbf{Y}_{t-1} | d, \mathcal{M}^*)} \left[ \log \frac{p(\mathbf{y}_t | \mathbf{Y}_{t-1}, d, \mathcal{M}^*)}{p(\mathbf{y}_t | \mathbf{Y}_{t-1}, d, \mathcal{M})} \right] \\ &\quad + \mathbb{E}_{p(\mathbf{Y}_{t-1} | d, \mathcal{M}^*)} \left[ \log \frac{p(\mathbf{Y}_{t-1} | d, \mathcal{M}^*)}{p(\mathbf{Y}_{t-1} | d, \mathcal{M})} \right] \\ &= \mathbb{E}_{p(\mathbf{Y}_{t-1} | d, \mathcal{M}^*)} [D_{\text{KL}}(p(\mathbf{y}_t | \mathbf{Y}_{t-1}, d, \mathcal{M}^*) || p(\mathbf{y}_t | \mathbf{Y}_{t-1}, d, \mathcal{M}))] \\ &\quad + \text{EDI}(d, \mathcal{M}, \mathcal{M}^*, t-1). \end{aligned} \quad (42)$$

Since the EDI is just a KL divergence, for the linear systems we have been studying in this paper, it is fairly straight forward to express it with the various quantities we have already derived. Therefore, we will state the main results without tenuous algebraic manipulation.

Single Step Update: For a single time step where data is generated by true process model  $p(\mathbf{y}_t | d, \mathbb{M}^*)$ , (see  $\mathbf{y}_t$  marginal of eq. 27), but we are evaluating  $\mathbb{M}$  according to

$p(\mathbf{y}_t | d, \mathbb{M})$  (see eq. 29), we can compute the KL divergence for these Gaussian distributions using eq. 5. Giving us:

$$\begin{aligned} \text{EDI}(d, \mathcal{M}, \mathcal{M}^*) &= \\ &\frac{1}{2} (\text{Tr}[\mathbf{S}_t^{-1} \mathbf{S}_t^*] + \log \frac{|\mathbf{S}_t|}{|\mathbf{S}_t^*|} + \Delta_y^T \mathbf{S}_t^{-1} \Delta_y - s). \end{aligned} \quad (43)$$

$s$  is the number of observations, e.g., sensors. Here we recall that  $\Delta_y = (\mathbf{H}^* \mu_{t|t-1}^* - \mathbf{H} \mu_{t|t-1})$  and emphasize that  $\mu_{t|t-1}^*$  and  $\mu_{t|t-1}$  need not be the same dimension since the comparison is happening in the data space.

For the special case were  $\mathbf{H}^* = [\mathbf{H}, \Delta]$ , the state of the model  $\mathcal{M}^*$  is  $\mathbf{x}_t^* = [\mathbf{x}_t, \delta_t]^T$ ,  $\mu_{\delta, t|t-1} = \mathbb{E}[\delta_{t|t-1}]$ , and  $\text{Cov}(\mathbf{x}_t^*) = \text{Diag}[\Sigma_{t|t-1}, \Gamma_{t|t-1}]$ , e.g., the augmented states are independent of other states. Then,

$$\begin{aligned} \text{EDI}(d, \mathcal{M}, \mathcal{M}^*) &= \frac{1}{2} (\text{Tr}[\mathbf{S}_t^{-1} \Delta \Gamma_{t|t-1} \Delta^T] \\ &- \log |\mathbb{I} + \mathbf{S}_t^{-1} \Delta \Gamma_{t|t-1} \Delta^T| + \mu_{\delta, t|t-1}^T \Delta^T \mathbf{S}_t^{-1} \Delta \mu_{\delta, t|t-1}) \end{aligned} \quad (44)$$

Infinite Horizon: For the asymptotic case, we may choose to ask a slightly different question when assessing the EDI. Instead of asking about a single  $\mathbf{y}_t$  we can ask about the full trajectory  $\mathbf{Y}_t = \{y_0 \dots y_t\}$ . Therefore to compute the EDI, we look to eq. 42. Under the previous assumptions of asymptotic stability, since we know that the predictives converge and are independent of the observations  $Y_t$ , we can expect the first term in eq. 42 to converge to a constant which we call  $\Delta_{\text{EDI}}$ . Therefore we expect  $\text{EDI}(d, \mathcal{M}, \mathcal{M}^*, t) \rightarrow t \Delta_{\text{EDI}}$  as  $t \rightarrow \infty$  unless  $\Delta_{\text{EDI}} = 0$ , meaning that there is only a finite amount of information to discriminate between the models based on the experiment even in the infinite horizon case. Therefore,  $\Delta_{\text{EDI}}$  is the critical quantity for understanding the asymptotic EDI. Using the expression for the Gaussian KL divergence, eq. 5 and taking the expectation using the asymptotic results found in Appendix VI-A, we find

$$\begin{aligned} \Delta_{\text{EDI}} &= \lim_{t \rightarrow \infty} \mathbb{E}_{p(\mathbf{Y}_{t-1} | d, \mathcal{M}^*)} [ \\ &D_{\text{KL}}(p(\mathbf{y}_t | \mathbf{Y}_{t-1}, d, \mathcal{M}^*) || p(\mathbf{y}_t | \mathbf{Y}_{t-1}, d, \mathcal{M}))] \\ &= \frac{1}{2} (\text{Tr}[\mathbf{S}^{-1} \mathbf{S}^*] + \log \frac{|\mathbf{S}|}{|\mathbf{S}^*|} + \text{Tr}[\mathbf{S}^{-1} M_S] - s). \end{aligned} \quad (45)$$

We recall that  $\mathbf{S}$  and  $\mathbf{S}^*$  are the stationary predictive covariances for an observation using design  $d$  for the models  $\mathcal{M}$  and  $\mathcal{M}^*$  respectively. The matrix  $M_S$  is given by,

$$M_S = \text{Cov}(\mathbf{H}^* \mu_{t|t-1}^* - \mathbf{H} \mu_{t|t-1})$$

$$= \begin{bmatrix} -HA & H^*A^* \end{bmatrix} M \begin{bmatrix} -HA & H^*A^* \end{bmatrix}^T \quad (46)$$

where  $M$  is the joint asymptotic covariance matrix of  $\mu_{t|t}$  and  $\mu_{t|t}^*$  and is the solution to the previously specified Lyapunov equation, eq. 37.

#### IV. EXAMPLES

##### A. Spring Mass Damper System

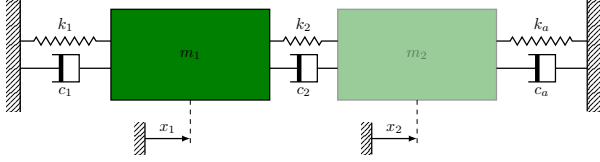


Fig. 1. Spring-Mass-Damper System with unknown second mass.

Fig. 1 shows a damped spring-mass system. The equations of motion for this system are,

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & m_1 & 0 \\ 0 & 0 & 0 & m_2 \end{bmatrix} \frac{d}{dt} \begin{bmatrix} x_1 \\ x_2 \\ v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ -(k_1 + k_2) & k_2 & -(b_1 + b_2) & b_2 \\ k_2 & -(k_2 + k_3) & b_2 & -(b_2 + b_3) \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ v_1 \\ v_2 \end{bmatrix} \quad (47)$$

where  $x_1, x_2$  denote the positions of the masses from their rest location. The variables  $v_1, v_2$  denote their linear velocities respectively. The spring constants are  $k_1, k_2, k_3$  and the damping coefficients are  $b_1, b_2, b_3$ . This continuous time linear system (CTLS) is then discretized for our analysis.

By analyzing the system we can see that under the conditions of high  $k_3$  stiffness, low  $m_2$  mass, or high  $b_3$  damping, that the two-mass system should behave close to a single-mass system. Therefore, under these conditions, we would expect the  $\Delta EDI$  criteria to become small when  $\mathcal{M}$  is the one-mass system and  $\mathcal{M}^*$  is the two-mass system. We see in panel A of Fig. 2 that  $\Delta EDI$  indeed decreases as we increase the stiffness  $k_3$ .

We now consider choosing an observer design  $d \in [0, \pi/2]$  to observe the position and velocity of the known mass,  $m_1$ , while balancing  $\Delta EDI$  and EIG. Our, admittedly arbitrary, observer measures the position and velocity of  $m_1$  with weights  $\cos(d)$  and  $\sin(d)$  respectively. The asymptotic EIG objective seeks to maximize information about the position and velocity of  $m_1$  according to  $\mathcal{M}$ . The  $\Delta EDI$  objective seeks to maximize our ability to asymptotically detect whether  $\mathcal{M}$  is plausible versus  $\mathcal{M}^*$ . Of course we don't know  $\mathcal{M}^*$  during the design phase so instead we average  $\Delta EDI$  over a prior range of stiffnesses from panel A of Fig. 2. We see how EIG varies over the designs as the navy-blue curve in Fig. 2, panel B, while the mean  $\Delta EDI$  is

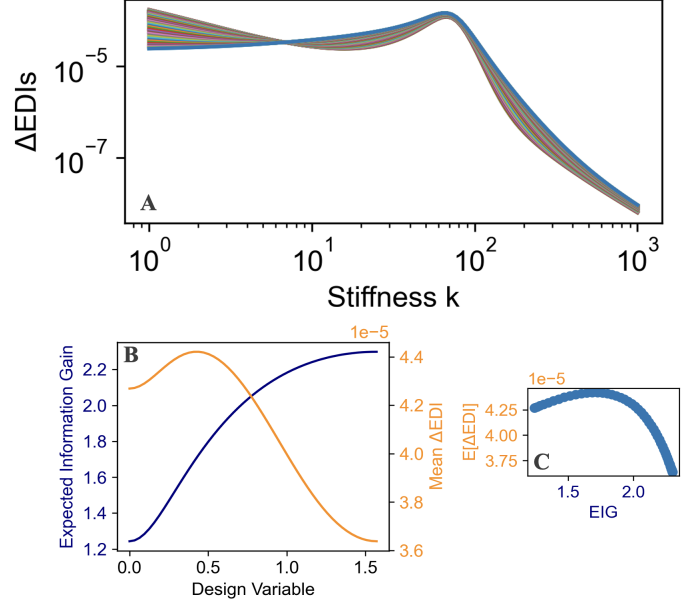


Fig. 2. Observer design and analysis for the spring mass system. Here, the true model  $\mathcal{M}^*$  is a two mass system while the inference model  $\mathcal{M}$  is the single mass system. Panel A shows how increasing the stiffness decreases our ability to distinguish between the models. Panels B and C show the trade off between EIG and  $\Delta EDI$  over our design variable.

shown as the orange curve. The trade off between these quantities is shown in panel C. Depending on the importance of discrimination vs performance, we may choose either to only observe the velocity (maximizing EIG) or to sacrifice some EIG to gain better discrimination power by choosing mixed sensor design.

##### B. F-16 Model

We use an F-16 aircraft model from [13] and [14]. This system originally has 12 states of which we pull out the

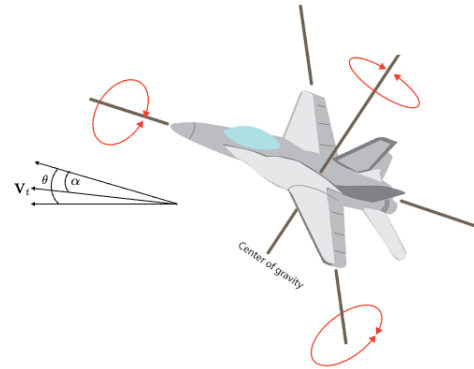


Fig. 3. F-16 model aircraft with specified states.

longitudinal dynamics with states (see Fig. 3):  $\theta$  (pitch angle),  $V$  (velocity),  $\alpha$  (attack angle),  $\dot{\theta}$  and controls:  $T$  (thrust),  $\delta_{ele}$  (elevator angle). We form a reduced-order CTLS using the closed loop system, which is then discretized.

For this model, we seek to add an additional output to the observer. This new output has the arbitrary form  $y_{new} = d_1\theta + d_2\alpha + d_3\dot{\theta}$ , where  $d_1^2 + d_2^2 + d_3^2 = 1$ . When considering these designs, we seek to balance maximizing asymptotic EIG while minimizing asymptotic EGIG. The inference model,  $\mathcal{M}$ , is the F-16 model with dynamics  $\mathbf{A}$ , but the true model,  $\mathcal{M}^*$ , has dynamics  $\mathbf{A}^* = \mathbf{A} + \Delta \odot \mathbf{A}$ . So,  $\mathbf{A}^*$  has perturbations scaled relative to  $\mathbf{A}$ . Because  $\Delta$  is unknown, we minimize the sensitivity of EGIG to changes of  $\Delta$ . Therefore, our metric is the norm,  $\|\nabla_{\Delta}\text{EGIG}(d_1, d_2)\|$ . The result is summarized in Fig. 4. Panel A shows the trade off of different designs between  $\text{EIG}(d_1, d_2)$  and  $\|\nabla_{\Delta}\text{EGIG}(d_1, d_2)\|$  and the Pareto front of optimal designs (purple). We see that the EGIG is much more sensitive to the design than the EIG, i.e., EGIG varies by about a factor of 4. Therefore, for a robust design we may sacrifice a little asymptotic EIG for meaningful improvement in robustness. Panels B and C show the EIG and the EGIG projected on the design space along with the corresponding Pareto set.

We have made the codes to these examples available on GitHub[15].

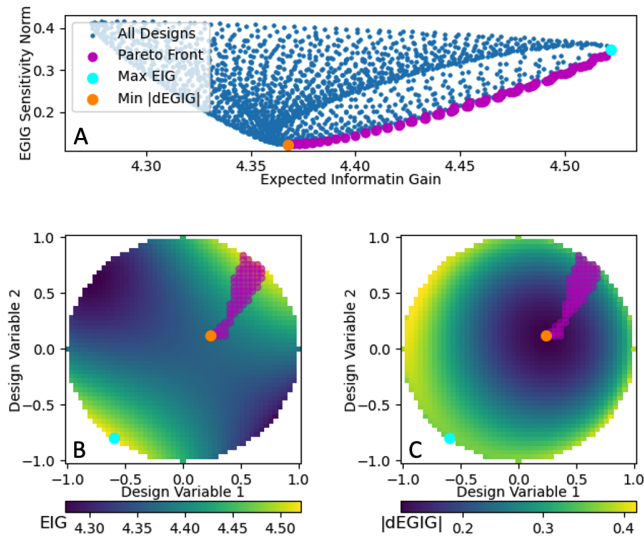


Fig. 4. Observer design for simplified F-16 model. The true model  $\mathcal{M}^*$  is a small perturbation,  $\Delta$ , in the dynamics from the inference model  $\mathcal{M}$ . We explore the addition of a new output,  $y_{new} = d_1\theta + d_2\alpha + d_3\dot{\theta}$  where  $d_3$  is constrained by  $d_1$  and  $d_2$ . We measure the improved performance using  $\text{EIG}(d_1, d_2)$  and robustness using the sensitivity of EGIG e.g.,  $\|\nabla_{\Delta}\text{EGIG}(d_1, d_2)\|$ . Panel A shows the trade off between these two criteria for different designs. Panels B and C show the projection of these criteria on to the design space.

## V. CONCLUSION

Maximizing the value of data for inference and prediction requires the careful selection of experimental conditions by modeling the experiment. These models are prone to misspecifications. We propose an information theoretic framework that extends the notion of Expected Information

Gain (EIG), typically used in Bayesian experiment design, to address the model mismatch issue. The proposed Expected Generalized Information Gain (EGIG) captures the information gain or loss with respect to a true model, when the experiment’s design is based on a model with discrepancy. On the other hand the proposed Expected Discriminatory Information (EDI) discriminates between models based upon the data generated, which further aids in model refinement. These three metrics are complementary as the EIG emphasizes data efficient experiments, the EGIG emphasizes experiments that lead to results that are robust to model discrepancy, and the EDI emphasizes experiments that would detect modeling failures.

For our future work we aim to develop a computational solver for assessing metrics within nonlinear systems and investigating the feasibility of its computational expense. Our objective is to establish correlations with alternative robustness measures like  $H_{\infty}$ , widely employed in control filter and observer design. Additionally, we aim to integrate the identification of worst-case scenarios, a departure from our current sensitivity metrics approach.

## ACKNOWLEDGMENT

This material is based upon work supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing conducted at Sandia National Laboratories. Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC., a wholly owned subsidiary of Honeywell International, Inc., for the U.S. Department of Energy’s National Nuclear Security Administration under contract DE-NA0003525. This paper describes objective technical results and analysis. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government. This research used resources of the National Energy Research Scientific Computing Center (NERSC), a U.S. Department of Energy Office of Science User Facility located at Lawrence Berkeley National Laboratory, operated under Contract No. DE-AC02-05CH11231 using NERSC award DOE-ERCAPm3876.

## REFERENCES

- [1] Jenny Brynjarsdottir and Anthony O’Hagan. “Learning about physical parameters: The importance of model discrepancy”. In: *Inverse problems* 30.11 (2014), p. 114007.
- [2] Peter Grünwald and Thijs van Ommen. “Inconsistency of Bayesian Inference for Misspecified Linear Models, and a Proposal for Repairing It”. In: *Bayesian Analysis* 12.4 (2017). DOI: 10.1214/17-ba1085.
- [3] Chi Feng. “Optimal Bayesian experimental design in the presence of model error”. MA thesis. Massachusetts Institute of Technology, 2015.

- [4] Sebastian Farquhar, Yarín Gal, and Tom Rainforth. *On Statistical Bias In Active Learning: How and When To Fix It*. 2021. arXiv: 2101.11665 [stat.ML].
- [5] Drew Fudenberg, Gleb Romanyuk, and Philipp Strack. “Active learning with a misspecified prior”. In: *Theoretical Economics* 12.3 (2017), pp. 1155–1189.
- [6] Jinwoo Go and Tobin Isaac. *Robust Expected Information Gain for Optimal Bayesian Experimental Design Using Ambiguity Sets*. 2022. arXiv: 2205.09914.
- [7] Antony M. Overstall and James M. McGree. *Bayesian decision-theoretic design of experiments under an alternative model*. 2021. arXiv: 1909.12570.
- [8] Sabina J. Sloman et al. *Characterizing the robustness of Bayesian adaptive experimental designs to active learning bias*. 2022. arXiv: 2205.13698.
- [9] Eric M Hernandez. “Balancing robustness and optimality in sensor placement for dynamic state estimation”. In: *Mechanical Systems and Signal Processing* 128 (2019), pp. 318–328.
- [10] Romain Pasquier and Ian FC Smith. “Robust system identification and model predictions in the presence of systematic uncertainty”. In: *Advanced Engineering Informatics* 29.4 (2015), pp. 1096–1109.
- [11] Jed A. Duersch and Thomas A. Catanach. “Generalizing Information to the Evolution of Rational Belief”. In: *Entropy* 22.1 (2020).
- [12] Josep Ginebra. “On the measure of the information in a statistical experiment”. In: (2007).
- [13] Brian L Stevens, Frank L Lewis, and Eric N Johnson. *Aircraft control and simulation: dynamics, controls design, and autonomous systems*. John Wiley & Sons, 2015.
- [14] Raktim Bhattacharya et al. “Nonlinear receding horizon control of an F-16 aircraft”. In: *Journal of Guidance, Control, and Dynamics* 25.5 (2002), pp. 924–931.
- [15] Niladri Das and Tommie A. Catanach. *Metrics BOED Model Misspecifications Python Notebooks*. Version 1.0.0. Mar. 2023. URL: <https://github.com/sandialabs/Metrics-BOED-Model-Misspecifications>.

## VI. APPENDIX

### A. Asymptotic Distribution of Inferred Means

Suppose we have a true model of a discrete time, asymptotically stable, linear dynamical system whose variables are denoted with a superscript  $*$ , while the model used for inference has variables without any superscript. Using a Kalman filter, the inferred mean is then given by,

$$\begin{aligned} \mu_t &= \mathbf{A}\mu_{t-1} + \mathbf{K}(Y_t - \mathbf{H}\mathbf{A}\mu_{t-1}) \\ &= (\mathbf{I} - \mathbf{K}\mathbf{H})\mathbf{A}\mu_{t-1} + \mathbf{K}\mathbf{H}^*\mathbf{A}^*\mu_{t-1}^* + \mathbf{K}\zeta^*. \end{aligned} \quad (48)$$

Here  $Y_t \sim \mathcal{N}(\mathbf{H}^*\mathbf{A}^*\mu_{t-1}^*, S^*)$ , so  $\zeta^* \sim \mathcal{N}(0, S^*)$ . Similarly, we can define the evolution of mean of the true

dynamical system under Kalman filtering as:

$$\begin{aligned} \mu_t^* &= \mathbf{A}^*\mu_{t-1}^* + \mathbf{K}^*(Y_t - \mathbf{H}^*\mathbf{A}^*\mu_{t-1}^*) \\ &= \mathbf{A}^*\mu_{t-1}^* + \mathbf{K}^*\zeta^* \end{aligned} \quad (49)$$

First, we note that  $E[\mu_t] = E[\mu_t^*] = 0$ , where the expectation is taken over asymptotically long sample trajectories of the true dynamical system. Second, we note that  $E[\mu_{t-1}\zeta^{*T}] = E[\mu_{t-1}^*\zeta^{*T}] = 0$ , i.e., they are independent. This comes from the fact that only  $\mu_{t-1}$  and  $\mu_{t-1}^*$  are functions of the trajectory and not  $\zeta^*$ , or in other words, all the information about the trajectory is captured in the mean estimates. Finally, we note that  $\mathbf{K}$  and  $\mathbf{K}^*$  are known for the asymptotic case by solving the respective DAREs, eq. 16-17. With that we can express the second moments of  $\mu_t$  and  $\mu_t^*$  as:

$$\mathbf{M}_t = \begin{pmatrix} E[\mu_t\mu_t^T] & E[\mu_t\mu_t^{*T}] \\ E[\mu_t^*\mu_t^T] & E[\mu_t^*\mu_t^{*T}] \end{pmatrix}. \quad (50)$$

In order to solve for the second moments, we define:

$$\mathbf{A} = \begin{pmatrix} (\mathbf{I} - \mathbf{K}\mathbf{H})\mathbf{A} & \mathbf{K}\mathbf{H}^*\mathbf{A}^* \\ \mathbf{0} & \mathbf{A}^* \end{pmatrix}, \quad (51)$$

$$\mathbf{Q} = \begin{pmatrix} \mathbf{K}\mathbf{S}^*\mathbf{K}^T & \mathbf{K}\mathbf{S}^*\mathbf{K}^{*T} \\ \mathbf{K}^*\mathbf{S}^*\mathbf{K}^T & \mathbf{K}^*\mathbf{S}^*\mathbf{K}^{*T} \end{pmatrix}. \quad (52)$$

Then we can solve for the moments as:

$$\mathbf{M}_t = \mathbf{A}\mathbf{M}_{t-1}\mathbf{A}^T + \mathbf{Q}. \quad (53)$$

Therefore, for the asymptotic case, we can solve the following Lyapunov equation to find the asymptotic second-order moments,

$$\mathbf{M} = \mathbf{A}\mathbf{M}\mathbf{A}^T + \mathbf{Q}, \quad (54)$$

giving us the result that asymptotically:

$$\begin{pmatrix} \mu_t \\ \mu_t^* \end{pmatrix} \sim \mathcal{N}(\mathbf{0}, \mathbf{M}). \quad (55)$$

*Special case  $\mathbf{A} = \mathbf{A}^*$ :* For the simpler case when we only have one model, the true model, we have the simplified equation given by:

$$\begin{aligned} \mathbf{M}^* &= \mathbf{A}^*\mathbf{M}^*\mathbf{A}^{*T} + \mathbf{K}^*\mathbf{S}^*\mathbf{K}^{*T} \\ &= \mathbf{A}^*\mathbf{M}^*\mathbf{A}^{*T} + \mathbf{A}^*\mathbf{P}_D^*\mathbf{A}^{*T} + \mathbf{Q} - \mathbf{P}_D^* \\ \Rightarrow \mathbf{M}^* + \mathbf{P}_D^* &= \mathbf{A}^*(\mathbf{M}^* + \mathbf{P}_D^*)\mathbf{A}^{*T} + \mathbf{Q}, \end{aligned} \quad (56)$$

the substitution  $\mathbf{K}^*\mathbf{S}^*\mathbf{K}^{*T} = \mathbf{A}^*\mathbf{P}_D^*\mathbf{A}^{*T} + \mathbf{Q} - \mathbf{P}_D^*$  can be found using the matrix inversion lemma and knowing that  $\mathbf{P}_D^*$  is the solution to eq. 12 for the asymptotic case. We observe that eq. 56 is a Lyapunov equation. Thus since we know that  $\mathbf{P}_L^*$  is the solution to the Lyapunov equation for this system. Therefore,

$$\mathbf{M}^* = \mathbb{E}[\mu_t^*\mu_t^{*T}] = \mathbf{P}_L^* - \mathbf{P}_D^*. \quad (57)$$