
Active Bayesian Causal Inference

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Causal discovery and causal reasoning are classically treated as separate and con-
2 secutive tasks: one first infers the causal graph, and then uses it to estimate causal
3 effects of interventions. However, such a two-stage approach is uneconomical, espe-
4 cially in terms of actively collected interventional data, since the causal query of inter-
5 est interest may not require a fully-specified causal model. From a Bayesian perspective,
6 it is also unnatural, since a causal query (e.g., the causal graph or some causal effect)
7 can be viewed as a latent quantity subject to posterior inference—quantities that are
8 not of direct interest ought to be marginalized out in this process, thus contributing
9 to our overall uncertainty. In this work, we propose Active Bayesian Causal Infer-
10 ence (ABCI), a *principled fully-Bayesian active learning framework for integrated*
11 *causal discovery and reasoning*, which jointly infers a posterior over causal models
12 and queries of interest. In our approach to ABCI, we focus on the class of causally-
13 sufficient nonlinear additive Gaussian noise models, which we model using Gaus-
14 sian processes. To capture the space of causal graphs, we use a continuous latent
15 graph representation, allowing our approach to scale to practically relevant problem
16 sizes. We sequentially design experiments that are maximally informative about our
17 target causal query, collect the corresponding interventional data, update our
18 beliefs, and repeat. Through simulations, we demonstrate that our approach is more
19 data-efficient than existing methods that only focus on learning the full causal graph.
20 This allows us to accurately learn downstream causal queries from fewer samples,
21 while providing well-calibrated uncertainty estimates of the quantities of interest.

22 1 Introduction

23 Causal reasoning, that is, answering causal queries such as the effect of a particular intervention, is
24 a fundamental scientific quest [3, 24, 27, 34]. A rigorous treatment of this quest requires a reference
25 causal model, typically consisting at least of (i) a causal diagram, or directed acyclic graph (DAG),
26 capturing the qualitative causal structure between a system’s variables [38]; and (ii) a joint distribution
27 which is Markovian w.r.t. this causal graph [52]. Other frameworks additionally model (iii) the func-
28 tional dependence of each variable on its causal parents in the graph [39, 58]. If the graph is not known
29 from domain expertise, *causal discovery* aims to infer it from data [33, 52]. However, given only obser-
30 vational (passively collected) data, causal discovery is fundamentally limited to recovering the Markov
31 equivalence class (MEC) of DAGs implying the same conditional independences as the data [52].
32 Additional structural assumptions (e.g., linearity) can render the graph identifiable [25, 42, 49, 59] but
33 are often hard to falsify, thus leading to risk of misspecification. These shortcomings motivate learning
34 from experimental (interventional) data which suffices to uniquely recover the true graph [10, 11, 19].
35 Here, we are particularly interested in the *active* learning setting in which we can sequentially design
36 and perform interventions that are most informative for the target causal query [1, 17, 19, 20, 35, 55].

37 Classically, causal discovery and reasoning are treated as separate, consecutive tasks that are studied
38 by different communities. Prior work on experimental design has thus focused either purely on causal
39 reasoning—how to best design experimental studies if the causal graph is known—or purely on causal
40 discovery, whenever the graph is unknown. In contrast, we consider the arguably more common
41 setting in which we are interested in performing causal reasoning but do not have access to a reference
42 causal model a priori. In this case, causal discovery can be seen as a means to an end, rather than as
43 the main objective. Nonetheless, existing experimental design approaches generally focus on learning

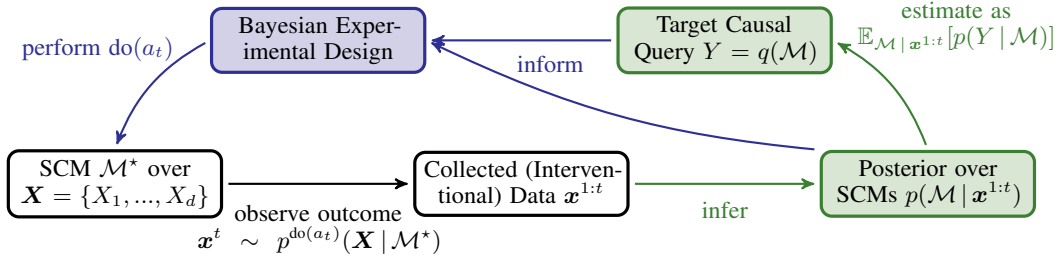


Figure 1: **Overview of the Active Bayesian Causal Inference (ABCI) framework.** At each time step t , we use Bayesian experimental design based on our current beliefs to choose a maximally informative intervention a_t to perform. We then collect a finite data sample from the interventional distribution induced by the environment, which we assume to be described by an unknown structural causal model (SCM) \mathcal{M}^* over a set of observable variables \mathbf{X} . Given (interventional) data $\mathbf{x}^{1:t}$ collected from the true SCM \mathcal{M}^* , together with a prior distribution over the model class of consideration, we infer the posterior over a target causal query $Y = q(\mathcal{M})$ that can be expressed as a function of the causal model: for example, we may be interested in the graph (causal discovery), the presence of certain edges (partial causal discovery), the full SCM (causal model learning), a collection of interventional distributions or treatment effects (causal reasoning), or any combination thereof.

44 the graph, which is subsequently fixed for the causal reasoning phase. This can be disadvantageous
 45 for two reasons: first, wasting samples on learning the full causal graph is suboptimal when we
 46 are only interested in specific aspects of the causal model; and second, causal discovery from finite
 47 (especially small amounts of) data entails significant epistemic model uncertainty—e.g., from low
 48 statistical test power or multiple highly-scoring DAGs—which should be taken into account [2, 13].

49 In the present work, we propose *Active Bayesian Causal Inference (ABCI)*, a principled, fully-
 50 Bayesian framework for integrated causal discovery and reasoning with experimental design. The
 51 basic approach is to put a Bayesian prior over the causal model class of choice, and to cast the
 52 learning problem as Bayesian inference over the model posterior. Moreover, we introduce the *target*
 53 *causal query* which is a function of the causal model that returns the (set of) causal quantities we are
 54 interested in. The model posterior together with the query function induce a *query posterior* which
 55 represents the result of our Bayesian learning procedure; it can be used, e.g., to derive a MAP solution
 56 or suitable expectation, or for down-stream decision tasks. The query posterior is incorporated in
 57 an active learning loop: we follow the Bayesian optimal experimental design approach [6, 28] and
 58 sequentially choose admissible interventions on the true causal model which are most informative
 59 about our target query w.r.t. our current beliefs. We then update our beliefs given the observed data by
 60 computing the posterior over causal models and queries, and use them to design the next experiment.

61 Since the general ABCI framework is computationally highly challenging, we implement it for the
 62 class of causally-sufficient nonlinear additive Gaussian noise models [25] which we model using
 63 Gaussian processes (GPs) [14, 57]. While this class is somewhat restrictive from a causal perspective,
 64 it is a flexible non-linear causal model which automates causal discovery in a wide range of scientific
 65 and engineering disciplines, as long as causal sufficiency can be reasonably assumed. To parameterize
 66 the combinatorial space of causal graphs, we use a recently proposed framework for differentiable
 67 Bayesian structure learning (DiBS) [30] that employs a continuous latent probabilistic graph
 68 representation to allow for tractable posterior inference. To approximately maximise information
 69 gain, we rely on Bayesian optimisation [31, 32, 51]. We highlight the following contributions:

- 70 • We propose ABCI as a flexible Bayesian active learning framework for efficiently inferring
 71 arbitrary sets of causal queries, subsuming causal discovery and reasoning as special cases (§ 3).
- 72 • We give a fully Bayesian treatment for the flexible class of nonlinear additive Gaussian noise
 73 models by leveraging GPs, continuous graph parametrisations, and Bayesian optimisation (§ 4).
- 74 • We demonstrate that our approach scales to relevant problem sizes and compares favourably to
 75 baselines in terms of efficiently learning the graph, full SCM, or interventional distributions (§ 5).

76 2 Related Work

77 Causal discovery and reasoning have been widely studied in machine learning and statistics [23, 42].
 78 Given an already collected set of observations, there is a large body of literature on learning causal
 79 structure, both in the form of a point estimate [18, 41, 49, 52] and a Bayesian posterior [2, 8, 13,
 80 21, 30]. Given a known causal graph, previous work studies how to estimate treatment effects
 81 or counterfactuals [39, 47, 48]. When interventional data is yet to be collected, existing work
 82 primarily focuses on the specific task of structure learning—without its downstream use. The concept

83 of (Bayesian) active causal discovery was first considered in discrete models with closed-form
 84 marginal likelihoods [35, 55] and later extended to nonlinear causal mechanisms [54, 56], multi-
 85 target interventions [53], and general models by using hypothesis testing [15]. Graph theoretic works
 86 give insights on the interventions required for full identifiability [10, 11, 19, 26].

87 Beyond learning the complete causal graph, few prior works have studied active causal inference.
 88 Concurrent work of Tigas et al. [54] considers experimental design for learning a full SCM
 89 parametrised by neural networks; there are significant differences to our approach: in particular, our
 90 framework (§ 3) is not limited to the information gain over the full model and provides a fully Bayesian
 91 treatment of the functions (§ 4). Agrawal et al. [1] consider actively learning a function of the causal
 92 graph under budget constraints, though not of the causal mechanisms and only for linear Gaussian
 93 models. Conversely, Rubenstein et al. [46] actively learn the causal mechanisms after the causal graph
 94 has been inferred. Thus, while prior work considers causal discovery and reasoning as a separate
 95 tasks, ABCI forms an integrated Bayesian approach for learning causal queries through interventions,
 96 reducing to previously studied settings in special cases. We further discuss related work in Appx. A.

97 3 Active Bayesian Causal Inference (ABCI) Framework

98 In this section, we first introduce the ABCI framework in generality, focusing on the main ideas
 99 and high-level ingredients, which are also illustrated in Fig. 1. In § 4 we then describe our particular
 100 implementation for the class of causally sufficient non-linear additive Gaussian noise models.

101 **Notation.** We use upper-case X and lower-case x to denote random variables and their realizations,
 102 respectively. Sets and vectors are written in bold face, \mathbf{X} and \mathbf{x} . With a slight abuse of notation, we
 103 use $p(\cdot)$ to denote different distributions, or densities, which are distinguished by their arguments.

104 **Causal Model.** To treat causality in a rigorous way, we first need to postulate a mathematically
 105 well-defined causal model. Historically hard questions about causality can then be reduced to
 106 *epistemic questions*, that is, what and how much is known about the causal model. A prominent
 107 type of causal model is the *structural causal model* (SCM) [39]. From a Bayesian perspective, an
 108 SCM can be viewed as a hierarchical data-generating process involving latent random variables.

109 **Definition 1** (SCM). An SCM \mathcal{M} over a set of *endogenous* (observed) variables $\mathbf{X} = \{X_1, \dots, X_d\}$
 110 and *exogenous* (latent) variables $\mathbf{U} = \{U_1, \dots, U_d\}$ consists of structural equations, or *mechanisms*,

$$X_i := f_i(\mathbf{Pa}_i, U_i), \quad \text{for } i \in \{1, \dots, d\}, \quad (3.1)$$

111 which assign the value of each X_i as a *deterministic* function f_i of its direct causes, or *causal parents*,
 112 $\mathbf{Pa}_i \subseteq \mathbf{X} \setminus \{X_i\}$ and U_i ; and a joint distribution $p(\mathbf{U})$ over the exogenous variables.

113 Associated with each SCM is a directed causal graph G with vertices \mathbf{X} and edges $X_j \rightarrow X_i$ iff.
 114 $X_j \in \mathbf{Pa}_i$, which we assume to be *acyclic* (i.e., it is a DAG). Any acyclic SCM then induces a
 115 unique *observational distribution* $p(\mathbf{X} | \mathcal{M})$ over the endogenous variables \mathbf{X} , which is obtained
 116 as the pushforward measure of $p(\mathbf{U})$ through the causal mechanisms in Eq. (3.1).

117 **Interventions.** A crucial aspect of causal models such as SCMs is that they also model the effect of
 118 *interventions*—external manipulations to one or more of the causal mechanisms in Eq. (3.1)—which,
 119 in general, are denoted using Pearl’s do-operator [39] as $\text{do}(\{X_i = \tilde{f}_i(\mathbf{Pa}_i, U_i)\}_{i \in \mathcal{I}})$ with $\mathcal{I} \subseteq [d]$
 120 and suitably chosen $\tilde{f}_i(\cdot)$. An intervention leads to a new SCM, the so-called *interventional SCM*,
 121 in which the relevant structural equations in Eq. (3.1) have been replaced by the new, manipulated
 122 ones. The interventional SCM thus induces a new distribution over the observed variables, the
 123 so-called *interventional distribution* which is denoted by $p^{\text{do}(a)}(\mathbf{X} | \mathcal{M})$ with a denoting the (set of)
 124 intervention(s) $\{X_i = \tilde{f}_i(\mathbf{Pa}_i, U_i)\}_{i \in \mathcal{I}}$. *Causal effects*—expressions like $\mathbb{E}[X_j | \text{do}(X_i = 3)]$ —can
 125 then be derived from the corresponding interventional distribution via standard probabilistic inference.

126 **Being Bayesian with Respect to Causal Models.** The main epistemic challenge for causal reasoning
 127 stems from the fact that the *true causal model* \mathcal{M}^* is not (or not completely) known. The canonical
 128 response to such epistemic challenges is a *Bayesian approach*: put a prior $p(\mathcal{M})$ on causal models,
 129 collect data \mathcal{D} from the true model \mathcal{M}^* , and compute the posterior via Bayes rule:

$$p(\mathcal{M} | \mathcal{D}) = \frac{p(\mathcal{D} | \mathcal{M}) p(\mathcal{M})}{p(\mathcal{D})} = \frac{p(\mathcal{D} | \mathcal{M}) p(\mathcal{M})}{\int p(\mathcal{D} | \mathcal{M}) p(\mathcal{M}) d\mathcal{M}}. \quad (3.2)$$

130 A full Bayesian treatment over \mathcal{M} is computationally delicate, to say the least. First, we require a
 131 way to parametrise the class of models \mathcal{M} we consider. Second, we need to be able to perform joint

132 posterior inference over this model class. In this paper, we present (one of) the first full Bayesian
 133 approaches which considers a flexible model class with nonlinear relationships (§ 4).

134 **Bayesian Causal Inference.** In the causal inference literature, the tasks of *causal discovery* (or,
 135 more generally, causal model learning) and *causal reasoning* are typically considered as separate
 136 problems. The former aims to learn (parts) of the causal model \mathcal{M}^* (typically the causal graph G^*)
 137 while the latter, assuming that the relevant parts of \mathcal{M}^* are already known, aims to identify and
 138 estimate some query of interest, typically using only observational data. This separation essentially
 139 suggests a two-stage approach: causal discovery followed by causal reasoning. From a Bayesian
 140 perspective, however, this distinction is unnatural and there is no real conceptual difference between
 141 the two. Rather, we might define a *causal query function* q , which specifies a *target causal query*
 142 $Y = q(\mathcal{M})$ as a function of the causal model \mathcal{M} . This view thus subsumes and generalises causal
 143 discovery and reasoning. Concretely, possible causal queries are

- 144 *Causal Discovery:* $Y = q_{\text{CD}}(\mathcal{M}) = G$, that is, learning the full causal graph G ;
- 145 *Partial Causal Discovery:* $Y = q_{\text{PCD}}(\mathcal{M}) = \phi(G)$, that is, learning some feature ϕ of the graph,
 146 such as the presence of a particular (set of) edge(s).
- 147 *Causal Model Learning:* $Y = q_{\text{CML}}(\mathcal{M}) = \mathcal{M}$, that is, learning the full SCM \mathcal{M} ;
- 148 *Causal Reasoning:* $Y = q_{\text{CR}}(\mathcal{M}) = \{p^{\text{do}(\mathbf{X}_{\mathcal{I}})}(X_j | \mathcal{M})\}_{j \in \mathcal{J}}$, that is, learning a set of
 149 interventional distributions induced by \mathcal{M} .¹

150 Once we have fixed the causal query, Bayesian inference naturally extends to the *query posterior*:

$$p(Y | \mathcal{D}) = \int p(Y | \mathcal{M}) p(\mathcal{M} | \mathcal{D}) d\mathcal{M} = \mathbb{E}_{\mathcal{M} | \mathcal{D}}[p(Y | \mathcal{M})], \quad (3.3)$$

151 where $p(Y | \mathcal{M})$ is deterministically given by $q(\mathcal{M})$, i.e., a point mass. Evidently, computing Eq. (3.3)
 152 constitutes a hard computational problem in general, as we need to marginalise over all causal
 153 models. In § 4 we introduce a practical implementation for a restricted causal model class.

154 **Identifiability of causal models and queries.** A crucial concept is that of *identifiability* of a model
 155 class, which refers to the ability to uniquely recover the true model in the limit of infinitely many
 156 observations from it [16].² In the context of our setting, if the class of causal models \mathcal{M} is identifiable,
 157 the model posterior $p(\mathcal{M} | \mathcal{D})$ in Eq. (3.2) and hence also the query posterior $p(Y | \mathcal{D})$ in Eq. (3.3)
 158 will collapse and converge to a point mass on their respective true values \mathcal{M}^* and $q(\mathcal{M}^*)$, given
 159 infinite data and provided the true model has non-zero mass under our prior, $p(\mathcal{M}^*) > 0$. Given
 160 only *observational* data, causal models are notoriously unidentifiable in general: without further
 161 assumptions on $p(\mathbf{U})$ and the structural form of Eq. (3.1), neither the graph nor the mechanisms can
 162 be recovered. In this case, $p(\mathcal{M} | \mathcal{D})$ may only converge to an equivalence class of models that cannot
 163 be further distinguished. Note, however, that even in this case, $p(Y | \mathcal{D})$ may still sometimes collapse,
 164 for example, if the Markov equivalence class (MEC) of graphs is identifiable (under causal sufficiency)
 165 and our query concerns the presence of a particular edge which is shared by all graphs in the MEC.

166 **Active Learning with Sequential Interventions.** Rather than collecting a *large, observational*
 167 dataset, we leverage observations from a *small* number of sequentially-performed *experiments*. The
 168 motivation for this is two-fold: first, experimental data can help resolve some of the non-identifiability
 169 issues discussed above; second, even if the model is identifiable (as for our approach in § 4), interven-
 170 tional data can still help learn our target causal query more quickly from *finite* data. Hence, at each
 171 time step t , we assume that we can perform an *experiment in the form of an intervention* a_t . The out-
 172 come of this experiment is a batch \mathbf{x}^t of N_t i.i.d. observations from the *true* interventional distribution:

$$\mathbf{x}^t = \{\mathbf{x}^{t,n}\}_{n=1}^{N_t}, \quad \mathbf{x}^{t,n} \stackrel{\text{i.i.d.}}{\sim} p^{\text{do}(a_t)}(\mathbf{X} | \mathcal{M}^*) \quad (3.4)$$

173 Note that restricting to $a_t = \emptyset$ —that is, sampling from the observational distribution—amounts
 174 to learning from observational data as a special case. Crucially, however, we design the experiment

¹Here the set \mathcal{J} can be uncountable, subsuming interventional distributions for a continuous set of interven-
 tions, possibly on different variables. Thus, in this case the return value of q is a set of density functions. In
 practice, these are implicitly represented in the learned Bayesian models, see § 5.

²It is worth pointing out that the term “identifiability” is sometimes used differently in the causal inference
 literature: within causal discovery, it typically refers to *structure identifiability*, that is, recovering only the causal
 graph; in the context of causal reasoning, on the other hand, it typically refers to whether an interventional (or
 counterfactual) query can be *expressed in terms of known quantities*, usually involving only the observational
 distribution. Here, we will use the term in its (original) statistical sense to refer to *identifiability of models*.

175 a_t so that it is *maximally informative* about our target causal query Y . In our Bayesian setting, this
 176 is naturally formulated by maximising the *information gain* between Y and the outcome \mathbf{X}^t [6, 28]:

$$\max_{a_t} \mathbb{I}(Y; \mathbf{X}^t | \mathbf{x}^{1:t-1}) \quad (3.5)$$

177 where \mathbf{X}^t follows the *predictive interventional distribution* of the Bayesian causal model ensemble
 178 at time $t - 1$ under intervention a_t , which is given by

$$\mathbf{X}^t \sim p^{\text{do}(a_t)}(\mathbf{X} | \mathbf{x}^{1:t-1}) \propto \int p^{\text{do}(a_t)}(\mathbf{X} | \mathcal{M}) p(\mathcal{M} | \mathbf{x}^{1:t-1}) d\mathcal{M}. \quad (3.6)$$

179 By maximizing Eq. (3.5) we collect experimental data in a goal-oriented manner to learn our causal
 180 query Y as efficiently and quickly as possible.

181 4 Tractable ABCI for Nonlinear Additive Noise Models

182 Having discussed the general framework and conceptual ideas, we now present our concrete approach
 183 to ABCI. This requires specifying: (i) the class of causal models we consider in Eq. (3.1), including
 184 their parametrisation; (ii) the types of interventions a_t we consider at each step and the corresponding
 185 interventional likelihood in Eq. (3.4); (iii) our prior distribution $p(\mathcal{M})$ over models; (iv) how to do
 186 posterior inference, that is, how to compute the model posterior in Eq. (3.2); and finally (v) how
 187 to maximise the information gain in Eq. (3.5) for experimental design.

188 **Model Class and Parametrisation.** In our approach to ABCI, we consider SCMs of the form

$$X_i := f_i(\mathbf{Pa}_i) + U_i, \quad \text{with} \quad U_i \sim \mathcal{N}(0, \sigma_i^2), \quad \text{for } i \in \{1, \dots, d\}, \quad (4.1)$$

189 where the f_i are *smooth, nonlinear* functions and where the U_i are assumed to be mutually
 190 independent, corresponding to the assumption of *causal sufficiency* (no hidden confounding). That
 191 is, we consider the special case of *causally sufficient, non-linear, Gaussian additive noise models*.
 192 Any model \mathcal{M} of this form can be described by a triple $\mathcal{M} = (G, \mathbf{f}, \boldsymbol{\sigma}^2)$, where G is a causal
 193 DAG, $\mathbf{f} = (f_1, \dots, f_d)$ is a vector of functions defined over the parent sets implied by G , and
 194 $\boldsymbol{\sigma}^2 = (\sigma_1^2, \dots, \sigma_d^2)$ contains the Gaussian noise variances. Provided that the f_i are nonlinear and
 195 not constant in any of their arguments, the model is identifiable almost surely [25, 43].

196 **Interventional Likelihood.** We support the realistic setting where only a subset $\mathbf{W} \subseteq \mathbf{X}$ of all
 197 variables are *actionable*, i.e., only \mathbf{W} can be the target of an intervention.³ For simplicity, we
 198 consider *hard interventions* of the form $\text{do}(a_t) = \text{do}(\mathbf{X}_{\mathcal{I}} = \mathbf{x}_{\mathcal{I}})$ which fix a subset $\mathbf{X}_{\mathcal{I}} \subseteq \mathbf{W}$ to a
 199 constant $\mathbf{x}_{\mathcal{I}}$. Due to causal sufficiency, the interventional likelihood under such hard interventions a_t
 200 factorises over the causal graph G and is given by the *g-formula* [44] or *truncated factorisation* [52]:

$$p^{\text{do}(a_t)}(\mathbf{X} | G, \mathbf{f}, \boldsymbol{\sigma}^2) = \mathbb{I}\{\mathbf{X}_{\mathcal{I}} = \mathbf{x}_{\mathcal{I}}\} \prod_{j \notin \mathcal{I}} p(X_j | f_j(\mathbf{Pa}_j^G), \sigma_j^2). \quad (4.2)$$

201 The last term in Eq. (4.2) is given by $\mathcal{N}(f_j(\mathbf{Pa}_j^G), \sigma_j^2)$ due to the Gaussian noise assumption.
 202 Let $\mathbf{x}^{1:t}$ be the entire dataset, collected up to time t . The likelihood of $\mathbf{x}^{1:t}$ is then given by

$$p(\mathbf{x}^{1:t} | G, \mathbf{f}, \boldsymbol{\sigma}^2) = \prod_{\tau=1}^t p^{\text{do}(a_\tau)}(\mathbf{x}^\tau | G, \mathbf{f}, \boldsymbol{\sigma}^2) = \prod_{\tau=1}^t \prod_{n=1}^{N_t} p^{\text{do}(a_\tau)}(\mathbf{x}^{\tau,n} | G, \mathbf{f}, \boldsymbol{\sigma}^2). \quad (4.3)$$

203 **Structured Model Prior.** To specify our model prior, we need to distinguish between *root nodes* X_i ,
 204 for which $\mathbf{Pa}_i = \emptyset$ and thus $f_i = \text{const}$, and *non-root nodes* X_j . For a given G , denote by
 205 $\mathbf{R}(G) = \{i \in [d] : \mathbf{Pa}_i^G = \emptyset\}$ the index set of root nodes, and by $\mathbf{NR}(G) = [d] \setminus \mathbf{R}(G)$ that of non-
 206 root nodes. We then place the following structured prior over the class of models $\mathcal{M} = (G, \mathbf{f}, \boldsymbol{\sigma}^2)$:

$$p(\mathcal{M}) = p(G) p(\mathbf{f}, \boldsymbol{\sigma}^2 | G) = p(G) \prod_{i \in \mathbf{R}(G)} p(f_i, \sigma_i^2 | G) \prod_{j \in \mathbf{NR}(G)} p(f_j | G) p(\sigma_j^2 | G). \quad (4.4)$$

207 Here, $p(G)$ is a prior over graphs, and $p(\mathbf{f}, \boldsymbol{\sigma}^2 | G)$ is a prior over the functions and noise variances
 208 in G . We factorise our prior conditional on G as in Eq. (4.4) not only to allow for a separate
 209 treatment of root nodes and non-root nodes but also to *share priors across similar graphs*: whenever
 210 $\mathbf{Pa}_i^{G_1} = \mathbf{Pa}_i^{G_2}$, we set $p(f_i, \sigma_i^2 | G_1) = p(f_i, \sigma_i^2 | G_2)$, and similarly for $p(f_j | G)$ and $p(\sigma_j^2 | G)$. As
 211 a consequence, the posteriors are also shared, which substantially reduces the computational burden.
 212 We also assume that $f_j \perp\!\!\!\perp f_{j'} | G$ and $\sigma_j^2 \perp\!\!\!\perp \sigma_{j'}^2 | G$ for all $j \neq j' \in \mathbf{NR}(G)$, which is motivated
 213 by the principle of independent causal mechanisms [42]. Our specific choices for $p(G)$, $p(f_i, \sigma_i^2 | G)$,
 214 $p(f_j | G)$, and $p(\sigma_j^2 | G)$ are guided by computational challenges and described in more detail below.

³In principle, the set of actionable variables might even change over time, in which case they are denoted \mathbf{W}_t .

215 **Model Posterior.** Given collected data $\mathbf{x}^{1:t}$, we can update our beliefs and quantify our uncertainty
 216 in \mathcal{M}^* by inferring a Bayesian posterior $p(\mathcal{M} | \mathbf{x}^{1:t})$ over SCMs $\mathcal{M} = (G, \mathbf{f}, \boldsymbol{\sigma}^2)$ as follows:⁴

$$p(\mathcal{M} | \mathbf{x}^{1:t}) = p(G | \mathbf{x}^{1:t}) \prod_{i \in \mathbf{R}(G)} p(f_i, \sigma_i^2 | \mathbf{x}^{1:t}, G) \prod_{j \in \mathbf{NR}(G)} p(f_j, \sigma_j^2 | \mathbf{x}^{1:t}, G). \quad (4.5)$$

217 For root nodes $i \in \mathbf{R}(G)$, posterior inference is straight-forward: we have $f_i = \text{const}$, so f_i can be
 218 viewed as the mean of U_i , cf. Eq. (4.1). We thus place a conjugate Normal-Gamma($\mu_i, \lambda_i, \alpha_i^R, \beta_i^R$)
 219 prior on $p(f_i, \sigma_i^2 | G)$, so that we can analytically compute the root node posterior $p(f_i, \sigma_i^2 | \mathbf{x}^{1:t}, G)$
 220 in Eq. (4.5) in closed form [36]. We collect all the Normal-Gamma hyperparameters in $(\boldsymbol{\mu}, \boldsymbol{\lambda}, \boldsymbol{\alpha}^R, \boldsymbol{\beta}^R)$.

221 The posteriors over graphs and non-root nodes $j \in \mathbf{NR}(G)$ are given as

$$p(G | \mathbf{x}^{1:t}) = \frac{p(\mathbf{x}^{1:t} | G) p(G)}{p(\mathbf{x}^{1:t})}, \quad p(f_j, \sigma_j^2 | \mathbf{x}^{1:t}, G) = \frac{p(\mathbf{x}^{1:t} | G, f_j, \sigma_j^2) p(f_j, \sigma_j^2 | G)}{p(\mathbf{x}^{1:t} | G)}. \quad (4.6)$$

222 Computing these posteriors is more involved and discussed below.

223 4.1 Addressing Challenges for Posterior Inference with GPs and DiBS

224 The particular challenges in Eq. (4.6) are the terms $p(\mathbf{x}^{1:t} | G)$ and $p(\mathbf{x}^{1:t})$. In the following, we
 225 will address these by means of appropriate prior choices and approximations.

226 **Challenge 1: Marginalising out Functions.** The term $p(\mathbf{x}^{1:t} | G)$ in Eq. (4.6) reads

$$p(\mathbf{x}^{1:t} | G) = \int p(\mathbf{x}^{1:t} | G, f_j, \sigma_j^2) p(f_j | G) p(\sigma_j^2 | G) df_j d\sigma_j^2 \quad (4.7)$$

227 and requires evaluating integrals over the function domain.
 228 We use *Gaussian processes* (GPs) [57] as an elegant
 229 way to solve this problem, as GPs can flexibly model
 230 *nonlinear* functions while offering convenient analytical
 231 properties. Specifically, we place a $\mathcal{GP}(0, k_j^G(\cdot, \cdot))$ prior
 232 on $p(f_j | G)$, where $k_j^G(\cdot, \cdot)$ is a covariance function over
 233 \mathbf{Pa}_j^G with length scales $\boldsymbol{\kappa}_j$, which we collect in $\boldsymbol{\kappa}$. In
 234 line with the GP-literature, we refer to $(\boldsymbol{\kappa}_j, \sigma_j^2)$ as the
 235 *GP-hyperparameters*. We place Gamma($\alpha_j^\sigma, \beta_j^\sigma$) and
 236 Gamma($\alpha_j^\kappa, \beta_j^\kappa$) priors on $p(\sigma_j^2 | G)$ and $p(\boldsymbol{\kappa}_j | G)$ and
 237 collect their parameters in $(\boldsymbol{\alpha}^{\text{GP}}, \boldsymbol{\beta}^{\text{GP}})$, see Fig. 2. For our
 238 model class, GPs then provide closed-form expressions
 239 for the ‘‘GP-marginal likelihood’’ $p(\mathbf{x}^{1:t} | G, \sigma_j^2, \boldsymbol{\kappa}_j)$, as
 240 well as for the ‘‘GP posteriors’’ $p(f_j | \mathbf{x}^{1:t}, G, \sigma_j^2, \boldsymbol{\kappa}_j)$, and
 241 the ‘‘predictive posteriors over observations’’ $p(\mathbf{X} | \mathbf{x}^{1:t}, G, \boldsymbol{\sigma}^2, \boldsymbol{\kappa})$ [57], see Appx. B for details.

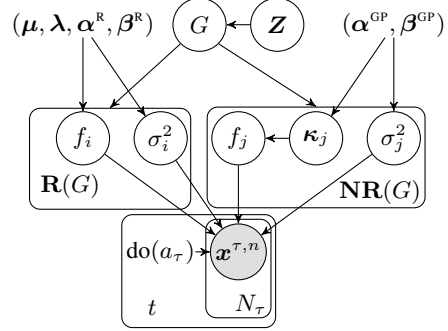


Figure 2: Graphical model representation of our GP-DiBS-ABCI approach.

242 **Challenge 2: Marginalising out GP-Hyperparameters.** While GPs allow for exact posterior inference
 243 conditional on a fixed value of $(\sigma_j^2, \boldsymbol{\kappa}_j)$, evaluating expressions such as $p(f_j | \mathbf{x}^{1:t}, G)$ requires
 244 marginalising out these GP-hyperparameters from the GP-posterior (see above). Unfortunately, this
 245 cannot, in general, be done exactly in connection with GPs as there is no closed-form expression for
 246 $p(\sigma_j^2, \boldsymbol{\kappa}_j | \mathbf{x}^{1:t}, G)$. We therefore approximate such expectations with a maximum a posteriori (MAP)
 247 point estimate $(\hat{\sigma}_j^2, \hat{\boldsymbol{\kappa}}_j)$, obtained by performing gradient ascent on the unnormalized log posterior,

$$\nabla \log p(\sigma_j^2, \boldsymbol{\kappa}_j | \mathbf{x}^{1:t}, G) = \nabla \log p(\mathbf{x}^{1:t} | G, \sigma_j^2, \boldsymbol{\kappa}_j) + \nabla \log p(\sigma_j^2, \boldsymbol{\kappa}_j | G) \quad (4.8)$$

248 according to a predefined update schedule, cf. Alg. 1. That is, we use approximations of the form:

$$p(f_j | \mathbf{x}^{1:t}, G) = \int p(f_j | \mathbf{x}^{1:t}, G, \sigma_j^2, \boldsymbol{\kappa}_j) p(\sigma_j^2, \boldsymbol{\kappa}_j | \mathbf{x}^{1:t}, G) d\sigma_j^2 d\boldsymbol{\kappa}_j \approx p(f_j | \mathbf{x}^{1:t}, G, \hat{\sigma}_j^2, \hat{\boldsymbol{\kappa}}_j)$$

249 **Challenge 3: Marginalising out Graphs.** Further, the ‘‘evidence’’ $p(\mathbf{x}^{1:t})$ is given by

$$p(\mathbf{x}^{1:t}) = \sum_G p(\mathbf{x}^{1:t} | G) p(G) \quad (4.9)$$

⁴To avoid further complicating the notation, we write all posteriors and likelihoods in terms of the full data $\mathbf{x}^{1:t}$. However, only observations of X_i and $X_j | \mathbf{Pa}_j^G$ matter for $i \in \mathbf{R}(G)$ and $j \in \mathbf{NR}(G)$.

Algorithm 1: GP-DiBS-ABCI for nonlinear additive Gaussian noise models

Input: no. of experiments T , batch sizes $\{N_t\}_{t=1}^T$, no. of latent particles M , no. of MC graphs K , particle resampling schedule $\{r_t\}_{t=1}^T$, hyperparameter update schedule $\{s_t\}_{t=1}^T$

Output: Posterior over target causal query $p(Y | \mathbf{x}^{1:T})$

for $t \leftarrow 1$ **to** T **do**

- $a_t \leftarrow \arg \max_{a=(\mathcal{I}, \mathbf{x}_{\mathcal{I}})} U(a, \mathbf{x}^{1:t-1})$ ▷design experiment: Eq. (4.11)
- $\mathbf{x}^t \leftarrow \{\mathbf{x}^{(t,n)} \sim p^{\text{do}(a_t)}(\mathbf{X} | \mathcal{M}^*)\}_{n=1}^{N_t}$ ▷perform experiment
- if** r_t **then**
 - $\mathbf{z}^t \leftarrow \text{resample_particles}(\mathbf{z}^t)$ ▷see App.D
- end**
- $\mathbf{G} \leftarrow \{\{G^{(k,m)} \sim p(G | \mathbf{z}_m)\}_{k=1}^K\}_{m=1}^M$ ▷sample graphs
- $\boldsymbol{\kappa}, \boldsymbol{\sigma}^2 \leftarrow \text{estimate_hyperparameters}(\mathbf{x}^{1:s_t}, \mathbf{G})$ ▷see Eq. (4.8)
- $\mathbf{z}^{t+1} \leftarrow \text{SVGD}(\mathbf{z}^t, \mathbf{x}^{1:t})$ ▷update latent particles

end

250 and involves a summation over all possible DAGs G . This becomes intractable even for $d \geq 4$
 251 variables as the number of DAGs grows super-exponentially in the number of variables [45]. To
 252 address this challenge, we employ the recently proposed DiBS framework [30]. By introducing
 253 a continuous prior $p(\mathbf{Z})$ that models G via $p(G | \mathbf{Z})$ and simultaneously enforces acyclicity of G ,
 254 Lorch et al. [30] show that we can efficiently infer the discrete posterior $p(G | \mathbf{x}^{1:t})$ via $p(\mathbf{Z} | \mathbf{x}^{1:t})$ as
 255

$$\mathbb{E}_{G | \mathbf{x}^{1:t}} [\phi(G)] = \mathbb{E}_{\mathbf{Z} | \mathbf{x}^{1:t}} \left[\frac{\mathbb{E}_{G | \mathbf{Z}} [p(\mathbf{x}^{1:t} | G) \phi(G)]}{\mathbb{E}_{G | \mathbf{Z}} [p(\mathbf{x}^{1:t} | G)]} \right] \quad (4.10)$$

256 where ϕ is some function of the graph. Since $p(\mathbf{Z} | \mathbf{x}^{1:t})$ is a continuous density with tractable
 257 gradient estimators, we can resort to efficient variational inference methods such as Stein Variational
 258 Gradient Descent (SVGD) for approximate inference [29], see Appx. D for additional details.

259 4.2 Approximate Bayesian Experimental Design with Bayesian Optimisation

260 As motivated in § 3, we aim to perform experiments a_t that are maximally informative about our
 261 target query $Y = q(\mathcal{M})$ by maximising the information gain from Eq. (3.5) given our current data
 262 $\mathcal{D} := \mathbf{x}^{1:t-1}$. In Appx. C we show that this is equivalent to maximising the following utility function:

$$U(a) = H(\mathbf{X}^t | \mathcal{D}) + \mathbb{E}_{\mathcal{M} | \mathcal{D}} \left[\mathbb{E}_{\mathbf{X}^t, Y | \mathcal{M}} \left[\log \mathbb{E}_{\mathcal{M}' | \mathcal{D}} \left[p(\mathbf{X}^t | \mathcal{M}') p(Y | \mathcal{M}') \right] \right] \right], \quad (4.11)$$

where $H(\mathbf{X}^t | \mathcal{D}) = \mathbb{E}_{\mathcal{M} | \mathcal{D}} \left[\mathbb{E}_{\mathbf{X}^t | \mathcal{M}} \left[\log \mathbb{E}_{\mathcal{M}' | \mathcal{D}} \left[p(\mathbf{X}^t | \mathcal{M}') \right] \right] \right]$

263 denotes the differential entropy of the experiment outcome, which depends on a and is distributed
 264 as in Eq. (3.6). This surrogate objective can be estimated using a nested Monte Carlo estimator, as
 265 long as we can sample from and compute $p(Y | \mathcal{M})$, see Appx. D for further details. For example,
 266 for $q_{\text{CR}}(\mathcal{M}) = p^{\text{do}(X_i=\psi)}(X_j | \mathcal{M})$ with $\psi \sim p(\psi)$ a distribution over intervention values, we get:

$$U_{\text{CR}}(a) = H(\mathbf{X}^t | \mathcal{D}) + \mathbb{E}_{\mathbf{X}^t | \mathcal{D}} \mathbb{E}_{\psi} \mathbb{E}_{X_j}^{\text{do}(X_i=\psi)} \left[\log \mathbb{E}_{\mathcal{M}' | \mathcal{D}} \left[p(\mathbf{X}^t | \mathcal{M}') p^{\text{do}(X_i=\psi)}(X_j | \mathcal{M}') \right] \right].$$

267 Importantly, for specific instances of the query function $q(\cdot)$ discussed in § 3, we can derive simpler
 268 utilities than Eq. (4.11). For example, for $q_{\text{CD}}(\mathcal{M}) = G$ and $q_{\text{CML}}(\mathcal{M}) = \mathcal{M}$ we arrive at

$$U_{\text{CD}}(a) = \mathbb{E}_{G | \mathcal{D}} \left[\mathbb{E}_{\mathbf{X}^t | G, \mathcal{D}} \left[\log p(\mathbf{X}^t | \mathcal{D}, G) - \log \mathbb{E}_{G' | \mathcal{D}} \left[p(\mathbf{X}^t | \mathcal{D}, G') \right] \right] \right], \quad (4.12)$$

$$U_{\text{CML}}(a) = \mathbb{E}_{\mathcal{M} | \mathcal{D}} \left[\mathbb{E}_{\mathbf{X}^t | \mathcal{M}} \left[\log p(\mathbf{X}^t | \mathcal{M}) - \log \mathbb{E}_{G' | \mathcal{D}} \left[p(\mathbf{X}^t | \mathcal{D}, G') \right] \right] \right], \quad (4.13)$$

269 where the entropy $\mathbb{E}_{\mathbf{X}^t | \mathcal{M}} [\log p(\mathbf{X}^t | \mathcal{M})]$ can again be efficiently computed given our modelling
 270 choices. For the sake of brevity, we defer derivations and estimation details to Appxs. C and D.

271 Finding the optimal experiment $a_t^* = (\mathcal{I}^*, \mathbf{x}_{\mathcal{I}^*}^*)$ requires jointly optimising the utility function cor-
 272 responding to our query with respect to (i) the set of intervention *targets* \mathcal{I} , and (ii) the corresponding
 273 intervention *values* $\mathbf{x}_{\mathcal{I}}$. This lends itself naturally to a nested, bi-level optimization scheme [56]:

$$\mathcal{I}^* \in \arg \max_{\mathcal{I}} U(\mathcal{I}, \mathbf{x}_{\mathcal{I}}^*), \quad \text{where} \quad \forall \mathcal{I} : \quad \mathbf{x}_{\mathcal{I}}^* \in \arg \max_{\mathbf{x}_{\mathcal{I}}} U(\mathcal{I}, \mathbf{x}_{\mathcal{I}}), \quad (4.14)$$

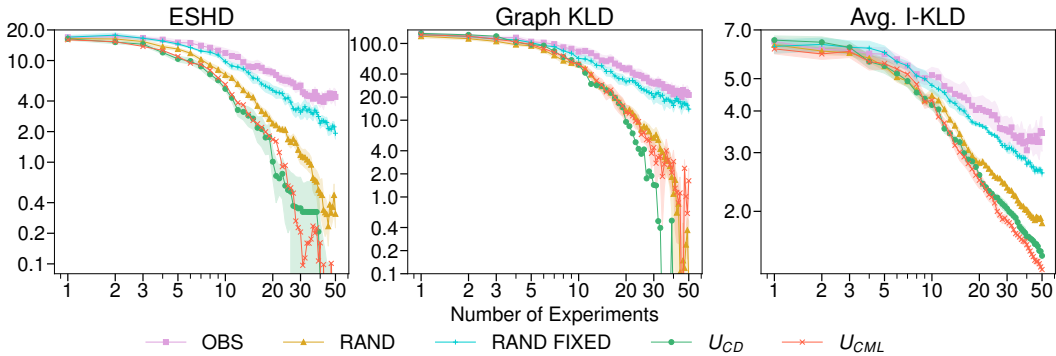


Figure 3: **Causal Discovery and SCM Learning.** Comparison of experimental design strategies for causal discovery (U_{CD}) and causal model learning (U_{CML}) with random and observational baselines on simulated ground truth models with 8 nodes. Lines and shaded areas show means ± 1 std. dev. across 30 runs (5 randomly sampled ground-truth SCMs with 6 restarts per SCM). (a) **ESHD.** Both our objectives significantly outperform the observational and random baselines. (b) **Graph-KLD.** U_{CD} , which optimises for this objective performs best as expected, but U_{CML} and the strong random baseline (RAND) perform competitively at learning the graph. (c) **Average I-KLD.** Both our strategies significantly outperform the baselines; U_{CML} , which aims to learn the full SCM, does slightly better than U_{CD} in terms of this proxy for causal model learning, as expected.

274 that is, we first estimate the optimal intervention values for all candidate intervention targets \mathcal{I} ,
 275 and then select the intervention target that yields the highest utility. The intervention target \mathcal{I}
 276 might contain multiple variables, which, however, yields a combinatorial problem. Thus, for
 277 simplicity, we consider only single-node interventions, i.e., $|\mathcal{I}| = 1$. To find $x_{\mathcal{I}}^*$, we employ *Bayesian*
 278 *optimisation* [31, 32, 51] to efficiently estimate an optimal intervention value $x_{\mathcal{I}}^*$, see Appx. D.

279 5 Experiments

280 **Setup.** We evaluate ABCI by inferring the query posterior on synthetic ground truth SCMs using
 281 several different experiment selection strategies. Specifically, we design experiments w.r.t. U_{CD}
 282 (causal discovery), U_{CML} (causal model learning), and U_{CR} (causal reasoning), see § 4.2. We compare
 283 against baselines which (i) only sample from the observational distribution (OBS) or (ii) pick an
 284 intervention target j uniformly at random from $[d] \cup \{\emptyset\}$ and set $X_j = 0$ (RAND FIXED, a weak
 285 random baseline used in prior work) or draw $X_j \sim \mathcal{U}(-7, 7)$ (RAND) if $X_j \neq \emptyset$. All methods
 286 follow our Bayesian GP-DiBS-ABCI approach from § 4. We sample ground truth SCMs over random
 287 scale-free graphs [4] of size $d = 8$, with mechanisms and noise variances drawn from our model
 288 prior Eq. (4.4). We initialise all methods with 5 observational samples, and then perform experiments
 289 with a batch size of 3. For specific prior choices and simulation details, see Appx. D.

290 **Metrics.** As ABCI infers a posterior over the target query Y , a natural evaluation choice is the
 291 Kullback-Leibler divergence (KLD) between the true query distribution and the inferred query poster-
 292 ior, $\text{KL}(p(Y | \mathcal{M}^*) || p(Y | \mathbf{x}^{1:t}))$. We report **Graph KLD**, a sample-based approximation of the
 293 KLD for posteriors over graphs (q_{CD}), and **Query KLD**, a KLD estimate for target interventional
 294 distributions (q_{CR}). As a proxy for the KLD of the SCM posterior (q_{CML}),⁵ we report the average
 295 KLD across all single node interventional distributions $\{p^{\text{do}(X_i=\psi)}(\mathbf{X})\}_{i=1}^d$, with $\psi \sim \mathcal{U}(-7, 7)$
 296 (**Average I-KLD**). We also report the *expected structural hamming distance* [9], **ESHD** =
 297 $\mathbb{E}_G | \mathbf{x}^{1:t} [\text{SHD}(G, G^*)]$, a commonly used causal discovery metric; see Appx. D for further details.

298 **Causal Discovery and SCM Learning (Fig. 3).** In our first experiment, we find that: (i) all our ABCI-
 299 based methods are able to meaningfully learn from small amounts of data, thus validating our Bayesian
 300 approach; further (ii) *performing targeted interventions using experimental design indeed yields*
 301 *improved performance over uninformed experimentation* (OBS, RAND FIXED, RAND). Notably, the
 302 stronger random baseline (RAND), which also randomises over intervention values, performs (surpris-
 303 ingly) well throughout—at least for the considered setting. As expected per the theoretical grounding
 304 of our information gain utilities, U_{CD} identifies the true graph the fastest (as measured by Graph
 305 KLD), whereas U_{CML} appears to most efficiently learn the full model, including the functions and
 306 noise variances, as measured by the Average I-KLD proxy, see the caption of Fig. 3 for further details.

307 **Learning Interventional Distributions (Fig. 4).** In our second experiment, we investigate ABCI’s
 308 causal reasoning capabilities by randomly sampling ground truth SCMs (as described above) over the
 309 fixed graph shown in Fig. 4 (right)—which is *not* known to the methods—and treat the (uncountable)

⁵The SCM KLD is either zero, if the SCM posterior collapses onto the true SCM, or infinite, otherwise.

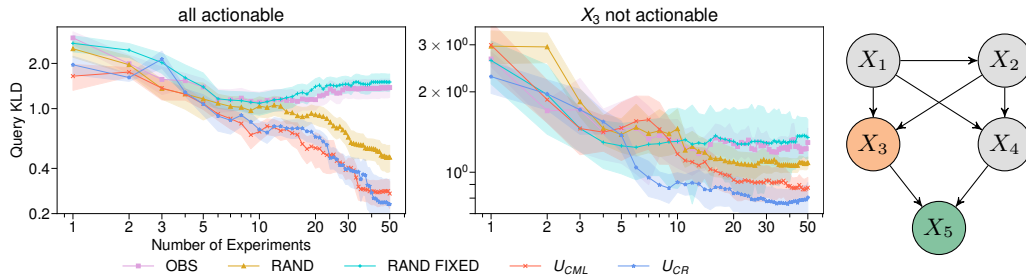


Figure 4: **Learning Interventional Distributions.** (left) Comparison of different methods w.r.t. learning a set of interventional distributions $p^{\text{do}(X_3=\psi)}(X_5 | \mathcal{M})$ with $\psi \sim \mathcal{U}[4, 7]$ on simulated ground truth models with fixed causal graph (right). Lines and shaded areas show mean ± 1 std. dev. across 25 runs (5 randomly sampled ground truth SCMs with 5 restarts each). **(a) All nodes actionable.** Our objectives significantly outperform the baselines; U_{CML} and U_{CR} perform similarly. In conjunction with results from Fig. 3, this suggests that U_{CML} yields a solid base model for performing downstream causal inference tasks. **(b) X_3 not actionable.** In this setting, where we cannot directly intervene on the treatment variable of interest, U_{CR} clearly outperforms all other methods for ≥ 5 experiments, suggesting that, in such a scenario, query-targeted experimental design is particularly helpful.

310 set of interventional distributions $p^{\text{do}(X_3=\psi)}(X_5 | \mathcal{M})$ with $\psi \sim \mathcal{U}[4, 7]$ as the target query. We find
 311 that *our informed experiment selection strategies significantly outperform the baselines at causal*
 312 *reasoning*, as measured by the Query KLD. In accord with the results from Fig. 3 and considering that
 313 that, once we know the true SCM, we can compute any causal quantity of interest, U_{CML} thus seems
 314 to provide a reasonable experimental strategy in case the causal query of interest is *not* known a
 315 priori. However, our results indicate that if we *do* know our query of interest, then U_{CR} *provides an*
 316 *even faster way for its estimation, especially when the treatment variable of interest is not directly*
 317 *intervenable*. Note the different axis scales, indicating that the task is harder in this case, as expected.

318 6 Discussion

319 **Assumptions, Limitations, and Extensions.** In § 4, we have made several assumptions to facilitate
 320 tractable inference and showcase the ABCI framework in a relatively simple causal setting. In
 321 particular, our assumptions exclude heteroscedastic noise, unobserved confounding, and cyclic
 322 relationships. On the experimental design side, we only considered *hard* interventions, but for
 323 some applications *soft* interventions [12] are more plausible. On the query side, we only considered
 324 *interventional* distributions. However, SCMs also naturally lend themselves to *counterfactual*
 325 reasoning, so one could also consider counterfactual queries such as the effect of the treatment
 326 on the treated [22, 50]. *In principle*, the ABCI framework as presented in § 3 extends directly to
 327 such generalisations. *In practice*, however, these can be non-trivial to implement, especially with
 328 regard to model parametrisation and tractable inference. Since actively performed interventions
 329 allow for causal learning even under causal sufficiency violations, we consider this a promising
 330 avenue for future work and believe the ABCI framework to be particularly well-suited for exploring it.
 331 Extensions to other causal modelling frameworks, such as graphical causal models are also of interest.

332 **Reflections on the ABCI Framework.** The main conceptual advantages of the ABCI framework
 333 are that it is *flexible* and *principled*. By considering general target causal queries, we can precisely
 334 specify what aspects of the causal model we are interested in, thereby offering a fresh perspective on
 335 the classical divide between causal discovery and reasoning: sometimes, the main objective may be
 336 to foster scientific understanding by uncovering the qualitative causal structure underlying real-world
 337 systems; other times, causal discovery may only be a means to an end—to support causal reasoning.
 338 Of particular interest in the context of actively selecting interventions is the setting where we cannot
 339 directly intervene on variables whose causal effect on others we are interested in (see Fig. 4), which
 340 connects to concepts such as transportability and external validity [5, 40]. ABCI is also flexible in
 341 that it easily allows for incorporating available domain knowledge: if we know some aspects of the
 342 model a priori (as assumed in conventional causal reasoning) or have access to a large observational
 343 sample (from which we can infer the MEC of DAGs), we can encode this in our prior and only
 344 optimise over a smaller model class, which should boost efficiency. The principled Bayesian nature
 345 of ABCI evidently comes at a significant computational cost: most integrals are intractable, and
 346 approximating them with Monte-Carlo sampling is computationally expensive and can introduce
 347 bias when resources are limited. On the other hand, in many real-world applications, such as in the
 348 context of biological networks, active interventions are possible but only at a significant cost [7, 37].
 349 Particularly in such cases, a careful and computationally-heavy experimental design approach as
 350 presented in the present work is warranted and might be easily amortised.

351 **References**

- 352 [1] Agrawal, R., Squires, C., Yang, K., Shanmugam, K., and Uhler, C. (2019). ABCD-strategy: Budgeted
353 experimental design for targeted causal structure discovery. In *The 22nd International Conference on Artificial*
354 *Intelligence and Statistics*, pages 3400–3409. PMLR.
- 355 [2] Agrawal, R., Uhler, C., and Broderick, T. (2018). Minimal I-MAP MCMC for scalable structure discovery
356 in causal DAG models. In *International Conference on Machine Learning*, pages 89–98. PMLR.
- 357 [3] Angrist, J. D. and Pischke, J.-S. (2008). *Mostly Harmless Econometrics*. Princeton University Press.
- 358 [4] Barabási, A.-L. and Albert, R. (1999). Emergence of scaling in random networks. *Science*, 286(5439):509–
359 512.
- 360 [5] Bareinboim, E. and Pearl, J. (2016). Causal inference and the data-fusion problem. *Proceedings of the*
361 *National Academy of Sciences*, 113(27):7345–7352.
- 362 [6] Chaloner, K. and Verdinelli, I. (1995). Bayesian experimental design: A review. *Statistical Science*, pages
363 273–304.
- 364 [7] Cho, H., Berger, B., and Peng, J. (2016). Reconstructing causal biological networks through active learning.
365 *PLoS one*, 11(3):e0150611.
- 366 [8] Cundy, C., Grover, A., and Ermon, S. (2021). BCD nets: Scalable variational approaches for Bayesian
367 causal discovery. *Advances in Neural Information Processing Systems*, 34.
- 368 [9] de Jongh, M. and Druzdzel, M. J. (2009). A comparison of structural distance measures for causal bayesian
369 network models. *Recent Advances in Intelligent Information Systems, Challenging Problems of Science,*
370 *Computer Science series*, pages 443–456.
- 371 [10] Eberhardt, F. (2008). Almost optimal intervention sets for causal discovery. In *Proceedings of the*
372 *Twenty-Fourth Conference on Uncertainty in Artificial Intelligence*, pages 161–168. AUAI Press.
- 373 [11] Eberhardt, F., Glymour, C., and Scheines, R. (2006). N-1 experiments suffice to determine the causal
374 relations among n variables. In *Innovations in machine learning*, pages 97–112. Springer.
- 375 [12] Eberhardt, F. and Scheines, R. (2007). Interventions and causal inference. *Philosophy of Science*,
376 74(5):981–995.
- 377 [13] Friedman, N. and Koller, D. (2003). Being Bayesian about network structure. a Bayesian approach to
378 structure discovery in Bayesian networks. *Machine learning*, 50(1):95–125.
- 379 [14] Friedman, N. and Nachman, I. (2000). Gaussian process networks. In *Proceedings of the Sixteenth*
380 *Conference on Uncertainty in Artificial Intelligence*, pages 211–219. Morgan Kaufmann Publishers Inc.
- 381 [15] Gamella, J. L. and Heinze-Deml, C. (2020). Active invariant causal prediction: Experiment selection
382 through stability. *Advances in Neural Information Processing Systems*, 33:15464–15475.
- 383 [16] George Casella, R. L. B. (2002). *Statistical Inference*, volume 2. Duxbury.
- 384 [17] Ghassami, A., Salehkaleybar, S., Kiyavash, N., and Bareinboim, E. (2017). Budgeted experiment design
385 for causal structure learning. *arXiv preprint arXiv:1709.03625*.
- 386 [18] Hauser, A. and Bühlmann, P. (2012). Characterization and greedy learning of interventional markov
387 equivalence classes of directed acyclic graphs. *The Journal of Machine Learning Research*, 13(1):2409–2464.
- 388 [19] Hauser, A. and Bühlmann, P. (2014). Two optimal strategies for active learning of causal models from
389 interventional data. *International Journal of Approximate Reasoning*, 55(4):926–939.
- 390 [20] He, Y.-B. and Geng, Z. (2008). Active learning of causal networks with intervention experiments and
391 optimal designs. *Journal of Machine Learning Research*, 9(Nov):2523–2547.
- 392 [21] Heckerman, D. (1995). A Bayesian approach to learning causal networks. In *Proceedings of the Eleventh*
393 *Conference on Uncertainty in Artificial Intelligence*, pages 285–295. Morgan Kaufmann Publishers Inc.
- 394 [22] Heckman, J. J. (1992). Policy evaluation. *Evaluating welfare and training programs*, page 201.
- 395 [23] Heinze-Deml, C., Maathuis, M. H., and Meinshausen, N. (2018). Causal structure learning. *Annual Review*
396 *of Statistics and Its Application*, 5:371–391.
- 397 [24] Hernán, M. A. and Robins, J. M. (2020). *Causal Inference: What If*. Boca Raton: Chapman & Hall/CRC.

- 398 [25] Hoyer, P. O., Janzing, D., Mooij, J. M., Peters, J., and Schölkopf, B. (2009). Nonlinear causal discovery
399 with additive noise models. In *Advances in neural information processing systems*, pages 689–696.
- 400 [26] Hyttinen, A., Eberhardt, F., and Hoyer, P. O. (2013). Experiment selection for causal discovery. *The*
401 *Journal of Machine Learning Research*, 14(1):3041–3071.
- 402 [27] Imbens, G. W. and Rubin, D. B. (2015). *Causal inference in statistics, social, and biomedical sciences*.
403 Cambridge University Press.
- 404 [28] Lindley, D. V. et al. (1956). On a measure of the information provided by an experiment. *The Annals of*
405 *Mathematical Statistics*, 27(4):986–1005.
- 406 [29] Liu, Q. and Wang, D. (2016). Stein variational gradient descent: A general purpose Bayesian inference
407 algorithm. In Lee, D., Sugiyama, M., Luxburg, U., Guyon, I., and Garnett, R., editors, *Advances in Neural*
408 *Information Processing Systems*, volume 29. Curran Associates, Inc.
- 409 [30] Lorch, L., Rothfuss, J., Schölkopf, B., and Krause, A. (2021). DiBS: Differentiable Bayesian Structure
410 Learning. *Advances in Neural Information Processing Systems*.
- 411 [31] Mockus, J. (1975). On Bayesian methods for seeking the extremum. In *Optimization Techniques IFIP*
412 *Technical Conference*, pages 400–404. Springer.
- 413 [32] Mockus, J. (2012). *Bayesian Approach to Global Optimization: Theory and Applications*, volume 37.
414 Springer Science & Business Media.
- 415 [33] Mooij, J. M., Peters, J., Janzing, D., Zscheischler, J., and Schölkopf, B. (2016). Distinguishing cause
416 from effect using observational data: methods and benchmarks. *The Journal of Machine Learning Research*,
417 17(1):1103–1204.
- 418 [34] Morgan, S. L. and Winship, C. (2014). *Counterfactuals and Causal Inference: Methods and Principles for*
419 *Social Research*. Cambridge University Press.
- 420 [35] Murphy, K. P. (2001). Active learning of causal Bayes net structure.
- 421 [36] Murphy, K. P. (2007). Conjugate Bayesian analysis of the gaussian distribution. Technical report, University
422 of British Columbia.
- 423 [37] Ness, R. O., Sachs, K., Mallick, P., and Vitek, O. (2017). A Bayesian active learning experimental design
424 for inferring signaling networks. In *International Conference on Research in Computational Molecular*
425 *Biology*, pages 134–156. Springer.
- 426 [38] Pearl, J. (1995). Causal diagrams for empirical research. *Biometrika*, 82(4):669–688.
- 427 [39] Pearl, J. (2009). *Causality*. Cambridge University Press, 2nd edition.
- 428 [40] Pearl, J. and Bareinboim, E. (2014). External validity: From do-calculus to transportability across
429 populations. *Statistical Science*, 29(4):579–595.
- 430 [41] Peters, J., Bühlmann, P., and Meinshausen, N. (2016). Causal inference by using invariant prediction:
431 identification and confidence intervals. *Journal of the Royal Statistical Society: Series B (Statistical*
432 *Methodology)*, 78(5):947–1012.
- 433 [42] Peters, J., Janzing, D., and Schölkopf, B. (2017). *Elements of Causal Inference - Foundations and Learning*
434 *Algorithms*. Adaptive Computation and Machine Learning Series. The MIT Press, Cambridge, MA, USA.
- 435 [43] Peters, J., Mooij, J. M., Janzing, D., and Schölkopf, B. (2014). Causal discovery with continuous additive
436 noise models. *The Journal of Machine Learning Research*, 15(1):2009–2053.
- 437 [44] Robins, J. (1986). A new approach to causal inference in mortality studies with a sustained exposure
438 period—application to control of the healthy worker survivor effect. *Mathematical modelling*, 7(9-12):1393–
439 1512.
- 440 [45] Robinson, R. W. (1973). Counting labeled acyclic digraphs. *New Directions in the Theory of Graphs*,
441 pages 239–273.
- 442 [46] Rubenstein, P. K., Tolstikhin, I., Hennig, P., and Schölkopf, B. (2017). Probabilistic active learning of
443 functions in structural causal models. *arXiv preprint arXiv:1706.10234*.
- 444 [47] Rubin, D. B. (2005). Causal inference using potential outcomes: Design, modeling, decisions. *Journal of*
445 *the American Statistical Association*, 100(469):322–331.

- 446 [48] Shalit, U., Johansson, F. D., and Sontag, D. (2017). Estimating individual treatment effect: generalization
447 bounds and algorithms. In *International Conference on Machine Learning*, pages 3076–3085. PMLR.
- 448 [49] Shimizu, S., Hoyer, P. O., Hyvärinen, A., Kerminen, A., and Jordan, M. (2006). A linear non-gaussian
449 acyclic model for causal discovery. *Journal of Machine Learning Research*, 7(10).
- 450 [50] Shpitser, I. and Pearl, J. (2009). Effects of treatment on the treated: Identification and generalization.
451 In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence, UAI 2009*, pages
452 514–521. AUAI Press.
- 453 [51] Snoek, J., Larochelle, H., and Adams, R. P. (2012). Practical Bayesian optimization of machine learning
454 algorithms. In *Advances in neural information processing systems*, pages 2951–2959.
- 455 [52] Spirtes, P., Glymour, C. N., and Scheines, R. (2000). *Causation, prediction, and search*. MIT press, 2nd
456 edition.
- 457 [53] Sussex, S., Uhler, C., and Krause, A. (2021). Near-optimal multi-perturbation experimental design for
458 causal structure learning. *Advances in Neural Information Processing Systems*, 34.
- 459 [54] Tigas, P., Annadani, Y., Jesson, A., Schölkopf, B., Gal, Y., and Bauer, S. (2022). Interventions, where and
460 how? experimental design for causal models at scale. *arXiv preprint arXiv:2203.02016*.
- 461 [55] Tong, S. and Koller, D. (2001). Active learning for structure in Bayesian networks. In *International Joint
462 Conference on Artificial Intelligence*, volume 17, pages 863–869.
- 463 [56] von Kügelgen, J., Rubenstein, P. K., Schölkopf, B., and Weller, A. (2019). Optimal experimental design
464 via Bayesian optimization: active causal structure learning for Gaussian process networks. In *NeurIPS
465 2019 Workshop “Do the right thing”: machine learning and causal inference for improved decision making*.
466 *arXiv:1910.03962*.
- 467 [57] Williams, C. K. and Rasmussen, C. E. (2006). *Gaussian Processes for Machine Learning*, volume 2. MIT
468 Press Cambridge, MA.
- 469 [58] Wright, S. (1934). The method of path coefficients. *The Annals of Mathematical Statistics*, 5(3):161–215.
- 470 [59] Zhang, K. and Hyvärinen, A. (2009). On the identifiability of the post-nonlinear causal model. In
471 *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, pages 647–655. AUAI
472 Press.

473 Checklist

- 474 1. For all authors...
- 475 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and
476 scope? [Yes] Summarizing the abstract/introduction we claim (i) to introduce a *principled fully-Bayesian*
477 *active learning framework for integrated causal discovery and reasoning* and to (ii) show the practicality
478 of our approach through simulations. We lay out the former concisely in § 3 and § 4. We report the
479 empirical evaluation in § 5.
- 480 (b) Did you describe the limitations of your work? [Yes] See discussion in § 6.
- 481 (c) Did you discuss any potential negative societal impacts of your work? [Yes] We provide a short
482 discussion in the Appendix.
- 483 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 484 2. If you are including theoretical results...
- 485 (a) Did you state the full set of assumptions of all theoretical results? [Yes] We give a concise and rigorous
486 treatment when formulating the general framework in § 3, as well as our approach and model specifics
487 in § 4.
- 488 (b) Did you include complete proofs of all theoretical results? [Yes] We provide the derivation of our utility
489 functions in Appx. C.
- 490 3. If you ran experiments...
- 491 (a) Did you include the code, data, and instructions needed to reproduce the main experimental results
492 (either in the supplemental material or as a URL)? [Yes] Python code and instructions are provided in
493 the supplement as source_code.zip
- 494 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
495 We give a minimal set of details in § 5 and provide full information about our experiments in Appx. D.
- 496 (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple
497 times)? [Yes] See Figs. 3 and 4

- 498 (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal
499 cluster, or cloud provider)? [Yes] We give a brief summary in Appendix D.
- 500 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 501 (a) If your work uses existing assets, did you cite the creators? [Yes] We do not use external models or data.
502 We use a set of Python packages that we list in Appendix D.
- 503 (b) Did you mention the license of the assets? [N/A]
- 504 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] We include our
505 code base in the supplementary material and will make it publicly available via Github upon acceptance.
- 506 (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating?
507 [N/A]
- 508 (e) Did you discuss whether the data you are using/curating contains personally identifiable information or
509 offensive content? [N/A]
- 510 5. If you used crowdsourcing or conducted research with human subjects...
- 511 (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
- 512 (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals,
513 if applicable? [N/A]
- 514 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant
515 compensation? [N/A]