

Manuscript version: Author's Accepted Manuscript

The version presented in WRAP is the author's accepted manuscript and may differ from the published version or Version of Record.

Persistent WRAP URL:

<http://wrap.warwick.ac.uk/142419>

How to cite:

Please refer to published version for the most recent bibliographic citation information.

Copyright and reuse:

The Warwick Research Archive Portal (WRAP) makes this work by researchers of the University of Warwick available open access under the following conditions.

Copyright © and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable the material made available in WRAP has been checked for eligibility before being made available.

Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

Publisher's statement:

Please refer to the repository item page, publisher's statement section, for further information.

For more information, please contact the WRAP Team at: wrap@warwick.ac.uk.

Multi-task Causal Learning with Gaussian Processes

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This paper studies the problem of learning the correlation structure of a set of
2 intervention functions defined on the directed acyclic graph (DAG) of a causal
3 model. This is useful when we are interested in jointly learning the causal effects of
4 interventions on different subsets of variables in a DAG, which is common in field
5 such as healthcare or operations research. We propose the first multi-task causal
6 Gaussian process (GP) model, which we call DAG-GP, that allows for information
7 sharing *across* continuous interventions and *across* experiments on different vari-
8 ables. DAG-GP accommodates different assumptions in terms of data availability
9 and captures the correlation between functions lying in input spaces of different
10 dimensionality via a well-defined integral operator. We give theoretical results
11 detailing *when* and *how* the DAG-GP model can be formulated depending on the
12 DAG. We test both the quality of its predictions and its calibrated uncertainties.
13 Compared to single-task models, DAG-GP achieves the best fitting performance in
14 a variety of real and synthetic settings. In addition, it helps to select optimal inter-
15 ventions faster than competing approaches when used within sequential decision
16 making frameworks, like active learning or Bayesian optimization.

17 1 Introduction

18 Solving decision making problems in a variety of domains such as healthcare, systems biology or
19 operations research, requires experimentation. By performing interventions one can understand
20 how a system behaves when an action is taken and thus infer the cause-effect relationships of a
21 phenomenon. For instance, in healthcare, drugs are tested in randomized clinical trials before
22 commercialization. Biologists might want to understand how genes interact in a cell once one of
23 them is knockout. Finally, engineers investigate the impact of design changes on complex physical
24 systems by conducting experiments on digital twins [33]. Experiments in these scenarios are usually
25 expensive, time-consuming, and, especially for field experiments, they may present ethical issues.
26 Therefore, researchers generally have to trade-off cost, time, and other practical considerations to
27 decide which experiments to conduct, if any, to learn about the system behaviour.

28 Consider the causal graph in Fig. 1 which describes how crop yield Y is affected by soil fumigants X
29 and the level of eel-worm population at different times $\mathbf{Z} = \{Z_1, Z_2, Z_3\}$ [11, 26]. By performing a
30 set of experiments, the investigator aims at learning the *intervention functions* relating the expected
31 crop yield to each possible intervention set and level. Naïvely, one could achieve that by modelling
32 each intervention function separately. However, this approach would disregard the correlation
33 structure existing across experimental outputs and would increase the computational complexity
34 of the problem. Indeed, the intervention functions are correlated and each experiment carries
35 information about the yield we would obtain by performing alternative interventions in the graph.
36 For instance, observing the yield when running an experiment on the *intervention set* $\{X, Z_1\}$ and
37 setting the value to the *intervention value* $\{x, z_1\}$, provides information about the yield we would
38 get from intervening only on X or on $\{X, Z_1, Z_2, Z_3\}$. This paper studies how to jointly model

such intervention functions so as to transfer knowledge across different experimental setups and integrate observational and interventional data. The model proposed here enables proper uncertainty quantification of the causal effects thus allowing to define optimal experimental design strategies.

1.1 Motivation and Contributions

The framework proposed in this work combines causal inference with multi-task learning via Gaussian processes (GP, [29]). Probabilistic causal models are commonly used in disciplines where explicit experimentation may be difficult and the *do*-calculus [26] allows to predict the effect of an intervention without performing the experiment. In *do*-calculus, different intervention functions are modelled individually and there is no information shared across experiments. Modelling the correlation across experiments is crucial especially when the number of observational data points is limited and experiments on some variables cannot be performed. Multi-task GP methods have been extensively used to model non-trivial correlations between outputs [4]. However, to the best of our knowledge, this is the first work focusing on intervention functions, possibly of different dimensionality, defined on a causal graph. Particularly, we make the following contributions:

- We give theoretical results detailing *when* and *how* a causal multi-task model for the experimental outputs can be developed depending on the topology of the DAG of a causal model.
- Exploiting our theoretical results, we develop a joint probabilistic model for all intervention functions, henceforth named DAG-GP, which flexibly accommodates different assumptions in terms of data availability – both observational and interventional.
- We demonstrate how DAG-GP achieves the best fitting performance in a variety of experimental settings while enabling proper uncertainty quantification and thus optimal decision making when used within Active Learning (AL) and Bayesian Optimization (BO).

1.2 Related work

While there exists an extensive literature on multi-task learning with GPs [9, 4] and causality [27, 17], the literature on causal multi-task learning is very limited. The majority of the studies have focused on domain adaptation problems [30, 25, 34] where data for a source domain is given, and the task is to predict the distribution of a target variable in a target domain. Several works [28, 6–8] have studied the problem of transferring the causal effects of a given variable *across* environments and have identified transportability conditions under which this is possible. Closer to our work, [2] have developed a linear coregionalization model for learning the individual treatment effects via observational data. While [2] is the first paper conceptualizing causal inference as a multi-task learning problem, its focus is on modelling the correlation across intervention levels for a single intervention function corresponding to a dichotomous intervention variable.

Differently from these previous works, this paper focuses on transfer *within* a single environment, *across* experiments and *across* intervention levels. The set of functions we wish to learn have continuous input spaces of different dimensionality. Therefore, capturing their correlation requires placing a probabilistic model over the inputs which enables mapping between input spaces. The DAG, which we assumed to be known and is not available in standard multi-task settings, allows us to define such a model. Therefore, *existing multi-output GP models are not applicable to our problem*.

Our work is also related to the literature on causal decision making. Studies in this field have focused on multi-armed bandit problems [5, 21, 24, 22] and reinforcement learning [10, 14] settings where arms or actions correspond to interventions on a DAG. More recently, [1] proposed a Causal Bayesian Optimization (CBO) framework solving the problem of finding an optimal intervention in a DAG by modelling the intervention functions with GPs. In CBO each function is modelled independently and their correlation is not accounted for when exploring the intervention space. This paper overcomes this limitation by introducing a multi-task model for experimental outputs. Finally, in the causal

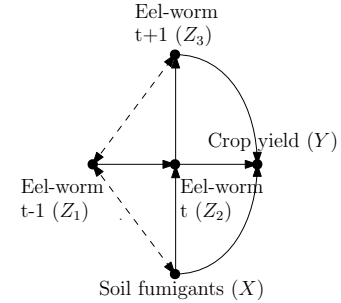


Figure 1: DAG for the crop yield. Nodes denote variables, arrows represent causal effects and dashed edges indicate unobserved confounders.

literature there has been a growing interest for experimental design algorithms to learn causal graphs [19, 18, 16] or the observational distributions in a graph [31]. Here we use our multi-task model within an AL framework so as to efficiently learn the experimental outputs in a causal graph.

2 Background and Problem setup

Consider a probabilistic structural causal model (SCM) [27] consisting of a directed acyclic graph \mathcal{G} (DAG) and a four-tuple $\langle \mathbf{U}, \mathbf{V}, F, P(\mathbf{U}) \rangle$, where \mathbf{U} is a set of independent *exogenous* background variables distributed according to the probability distribution $P(\mathbf{U})$, \mathbf{V} is a set of observed *endogenous* variables and $F = \{f_1, \dots, f_{|\mathbf{V}|}\}$ is a set of functions such that $v_i = f_i(\text{Pa}_i, u_i)$ with $\text{Pa}_i = \text{Pa}(V_i)$ denoting the parents of V_i . \mathcal{G} encodes our knowledge of the existing causal mechanisms among \mathbf{V} . Within \mathbf{V} , we distinguish between two different types of variables: treatment variables \mathbf{X} that can be manipulated and set to specific values¹ and output variables \mathbf{Y} that represent the agent’s outcomes of interest. Given \mathcal{G} , we denote the *interventional distribution* for two disjoint sets in \mathbf{V} , say \mathbf{X} and \mathbf{Y} , as $P(\mathbf{Y}|\text{do}(\mathbf{X} = \mathbf{x}))$. This is the distribution of \mathbf{Y} obtained by intervening on \mathbf{X} and fixing its value to \mathbf{x} in the data generating mechanism, irrespective of the values of its parents. The interventional distribution differs from the *observational distribution* which is denoted by $P(\mathbf{Y}|\mathbf{X} = \mathbf{x})$. Under some identifiability conditions [15], *do*-calculus allows to estimate interventional distributions and thus causal effects from observational distributions [26]. In this paper, we assume the causal effect for \mathbf{X} on \mathbf{Y} to be identifiable $\forall \mathbf{X} \in \mathcal{P}(\mathbf{X})$ with $\mathcal{P}(\mathbf{X})$ denoting the power set of \mathbf{X} .

2.1 Problem setup

Consider a DAG \mathcal{G} and the related SCM. Define the set of intervention functions for Y in \mathcal{G} as:

$$\mathbf{T} = \{t_s(\mathbf{x})\}_{s=1}^{|\mathcal{P}(\mathbf{X})|} \quad t_s(\mathbf{x}) = \mathbb{E}_{p(Y|\text{do}(\mathbf{X}_s = \mathbf{x}))}[Y] = \mathbb{E}[Y|\text{do}(\mathbf{X}_s = \mathbf{x})]. \quad (1)$$

with $\mathbf{X}_s \in \mathcal{P}(\mathbf{X})$ where $\mathcal{P}(\mathbf{X})$ is the power set of \mathbf{X} minus the empty set² and $\mathbf{x} \in D(\mathbf{X}_s)$ where $D(\mathbf{X}_s) = \times_{X \in \mathbf{X}_s} D(X)$ with $D(X)$ denoting the *interventional domain* of X . Let $\mathcal{D}^O = \{\mathbf{x}_n, y_n\}_{n=1}^N$, with $\mathbf{x}_n \in \mathbb{R}^{|\mathbf{X}|}$ and $y_n \in \mathbb{R}$, be an observational dataset of size N from this SCM.

Consider an interventional dataset $\mathcal{D}^I = (\mathbf{X}^I, \mathbf{Y}^I)$ with $\mathbf{X}^I = \bigcup_s \{\mathbf{x}_{si}^I\}_{i=1}^{N_s^I}$ and $\mathbf{Y}^I = \bigcup_s \{y_{si}^I\}_{i=1}^{N_s^I}$ denoting the intervention levels and the function values observed from previously run experiments across sets in $\mathcal{P}(\mathbf{X})$. N_s^I represents the number of experimental outputs observed for the intervention set \mathbf{X}_s . Our goal is to define a joint prior distribution $p(\mathbf{T})$ and compute the posterior $p(\mathbf{T}|\mathcal{D}^I)$ so as to make probabilistic predictions for \mathbf{T} at some unobserved intervention sets and levels.

3 Multi-task learning of intervention functions

In this section we address the following question: *can we develop a joint model for the functions \mathbf{T} in a causal graph and thus transfer information across experiments?*

To answer this question we study the correlation among functions in \mathbf{T} which varies with the topology of \mathcal{G} . Inspired by previous works on latent force models [3], we show how any functions in \mathbf{T} can be written as an integral transformation of some base function f , also defined starting from \mathcal{G} , via some integral operator L_s such that $t_s(\mathbf{x}) = L_s(f)(\mathbf{x})$, $\forall \mathbf{X}_s \in \mathcal{P}(\mathbf{X})$. We first characterize the latent structure among experimental outputs and provide an explicit expression for both f and L_s for each intervention set (§3.1). Based on the properties of \mathcal{G} , we clarify when this function exists. Exploiting these results, we detail a new model to learn \mathbf{T} which we call the DAG-GP model (§3.2). In DAG-GP we place a GP prior on f and propagate our prior assumptions on the remaining part of the graph to analytically derive a joint distribution of the elements in \mathbf{T} . The resulting prior distribution incorporates the causal structure and enables the integration of observational and interventional data.

3.1 Characterization of the latent structure in a DAG

Next results provide a theoretical foundation for the multi-task causal GP model introduced later. In particular, they characterize when f and L_s exist and how to compute them thus fully characterizing when transfer across experiments is possible. All proofs are given in the appendix.

¹This setting can be extended to include non-manipulative variables. See [23] for a definition of such nodes.

²We exclude the empty set as it corresponds to the observational distribution $t_\emptyset(\mathbf{x}) = \mathbb{E}[Y]$.

136 **Definition 3.1.** Consider a DAG \mathcal{G} where the treatment variables are denoted by \mathbf{X} . Let \mathbf{C} be the set
 137 of variables directly confounded with Y , \mathbf{C}^N be the set of variables in \mathbf{C} that are not colliders³ and \mathbf{I}
 138 be the set $\text{Pa}(Y)$. For each $\mathbf{X}_s \in \mathcal{P}(\mathbf{X})$ we define the following sets:

- 139 • $\mathbf{I}_s^N = \mathbf{I} \setminus (\mathbf{X}_s \cap \mathbf{I})$ represents the set of variables in \mathbf{I} not included in \mathbf{X}_s .
- 140 • $\mathbf{C}_s^I = \mathbf{C}^N \cap \mathbf{X}_s$ is the set of variables in \mathbf{C} which are included in \mathbf{X}_s and are not colliders.
- 141 • $\mathbf{C}_s^N = \mathbf{C}^N \setminus \mathbf{C}_s^I$ is the set of variables in \mathbf{C} that are neither included in \mathbf{X}_s nor colliders.

142 In the following theorem \mathbf{v}_s^N gives the values for the variables in the set \mathbf{I}_s^N while \mathbf{c} represents the
 143 values for the set \mathbf{C}^N which are partition in \mathbf{c}_s^N and \mathbf{c}_s^I depending on the set \mathbf{X}_s we are considering.

144 **Theorem 3.1. Causal operator.** Consider a causal graph \mathcal{G} and a related SCM where the output
 145 variable and the treatment variables are denoted by Y and \mathbf{X} respectively. Denote by \mathbf{C} the set
 146 of variables in \mathcal{G} that are directly confounded with Y and let \mathbf{I} be the set $\text{Pa}(Y)$. Assume that \mathbf{C}
 147 does not include nodes that have both unconfounded incoming and outgoing edges. It is possible
 148 to prove that, $\forall \mathbf{X}_s \in \mathcal{P}(\mathbf{X})$, the intervention function $t_s(\mathbf{x}) : D(\mathbf{X}_s) \rightarrow \mathbb{R}$ can be written as
 149 $t_s(\mathbf{x}) = L_s(f)(\mathbf{x})$ where

$$L_s(f)(\mathbf{x}) = \int \cdots \int \pi_s(\mathbf{x}, (\mathbf{v}_s^N, \mathbf{c})) f(\mathbf{v}, \mathbf{c}) d\mathbf{v}_s^N d\mathbf{c}, \quad (2)$$

150 with $f(\mathbf{v}, \mathbf{c}) = \mathbb{E}[Y | do(\mathbf{I} = \mathbf{v}), \mathbf{C}^N = \mathbf{c}]$ representing a shared latent function and
 151 $\pi_s(\mathbf{x}, (\mathbf{v}_s^N, \mathbf{c})) = p(\mathbf{c}_s^I | \mathbf{c}_s^N) p(\mathbf{v}_s^N, \mathbf{c}_s^I | do(\mathbf{X}_s = \mathbf{x}))$ giving the integrating measure for the set \mathbf{X}_s .

152 In the sequel we call $L_s(f)(\mathbf{x})$ the *causal operator*, $(\mathbf{I} \cup \mathbf{C})$ the *base set*, $f(\mathbf{v}, \mathbf{c})$ the *base function*
 153 and $\pi_s(\cdot, \cdot)$ the *integrating measure* of the set \mathbf{X}_s . A simple limiting case arises when the DAG does
 154 not include variables directly confounded with Y or \mathbf{C} only includes colliders. In this case $\mathbf{C} = \emptyset$
 155 and the base function is included in \mathbf{T} . Theorem 3.1 provides a mechanism to reconstruct all causal
 156 effects emerging from $\mathcal{P}(\mathbf{X})$ using the base function as a “driving force”. In particular, the integrating
 157 measures can be seen as Green’s functions incorporating the DAG structure [3]. While it can be
 158 further generalized to select \mathbf{I} to be different from $\text{Pa}(Y)$, this choice is particularly useful due to the
 159 following result.

160 **Corollary 3.1. Minimality of \mathbf{I} .** The smallest set \mathbf{I} for which Eq. (2) holds is given by $\text{Pa}(Y)$.

161 The dimensionality of \mathbf{I} when chosen as $\text{Pa}(Y)$ has properties that have been previously studied
 162 in the literature. In the context of optimization [1], it corresponds to the so-called causal intrinsic
 163 dimensionality, which refers to the effective dimensionality of the space in which a function is
 164 optimized when causal information is available. The existence of f depends on the properties of the
 165 nodes in \mathbf{C} which also represents the smallest set for which Eq. (2) holds (§1.4 in the supplement).

166 **Theorem 3.2. Existence of f .** If \mathbf{C} includes nodes that have both unconfounded incoming and
 167 outgoing edges the function f does not exist.

168 When f does not exist, full transfer across all functions in \mathbf{T} is not possible (DAGs with red edges in
 169 Fig. 4). However, these results enable a model for partial transfer across a subset of \mathbf{T} (§1.6 supp.).

170 3.2 The DAG-GP model

171 Next, we introduce the DAG GP model based on the results from the previous section.

172 **Model Likelihood:** Let $\mathcal{D}^I = (\mathbf{X}^I, \mathbf{Y}^I)$ be the interventional dataset defined in Section 2.1. Denote
 173 by \mathbf{T}^I the collection of intervention vector-valued functions computed at \mathbf{X}^I . Each entry y_{si}^I in \mathbf{Y}^I ,
 174 is assumed to be a noisy observation of the corresponding function t_s at \mathbf{x}_i^I :

$$y_{si}^I = t_s(\mathbf{x}_i^I) + \epsilon_{si}, \text{ for } s = 1, \dots, |\mathcal{P}(\mathbf{X})| \text{ and } i = 1, \dots, N_s^I, \quad (3)$$

175 with $\epsilon_{si} \sim \mathcal{N}(0, \sigma^2)$. In compact form, the joint likelihood function is $p(\mathbf{Y}^I | \mathbf{T}^I, \sigma^2) = \mathcal{N}(\mathbf{T}^I, \sigma^2 \mathbf{I})$.

176 **Prior distribution on \mathbf{T} :** To define a join prior on the set of intervention functions, $p(\mathbf{T})$, we take
 177 the following steps. First, we follow [1] to place a *causal prior* on f , the base function of the DAG.
 178 Second, we propagate this prior on f through all elements in \mathbf{T} via the causal operator in Eq. (2).

³Variables in \mathbf{C} causally influenced by \mathbf{X} and Y .

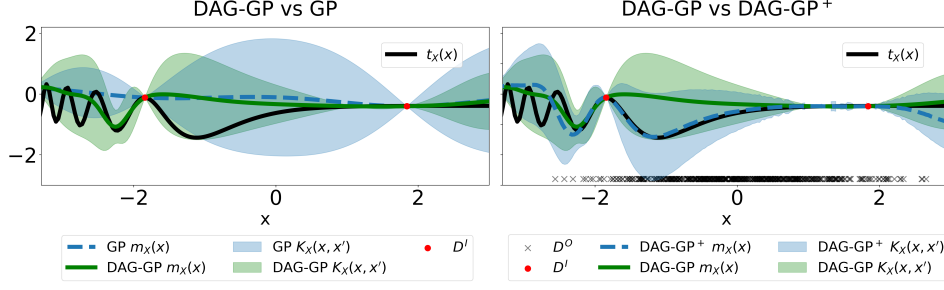


Figure 2: Posterior mean and variance for $t_X(\mathbf{x})$ in the DAG of Fig. 4 (a) (without the red edge). For both plots $m_X(\cdot)$ and $K_X(\cdot, \cdot)$ give the posterior mean and standard deviation respectively. *Left:* Comparison between the DAG-GP model and a single-task GP model (GP). DAG-GP captures the behaviour of $t_X(\mathbf{x})$ in areas where \mathcal{D}^I is not available (see area around $x = -2$) while reducing the uncertainty via transfer due to available data for \mathbf{z} . *Right:* Comparison between DAG-GP with the causal prior (DAG-GP⁺) and a standard prior with zero mean and RBF kernel (DAG-GP). In addition to transfer, DAG-GP⁺ captures the behaviour of $t_X(\mathbf{x})$ in areas where \mathcal{D}^O (black \times) is available (see region $[-2, 0]$) while inflating the uncertainty in areas with no observational data.

179 *Step 1, causal prior on the base function:* The key idea of the causal prior, already used in [1], is to
 180 use the observational dataset \mathcal{D}^O and the *do*-calculus to construct the prior mean and variance of a
 181 GP that is used to model an intervention function. Our aim is to compute such prior for the causal
 182 effect of the base set $\mathbf{I} \cup \mathbf{C}$ on Y . The causal prior has the benefit of carrying causal information but
 183 at the expense of requiring \mathcal{D}^O to estimate the causal effect. Any sensible prior can be used in this
 184 step, so the availability of \mathcal{D}^O is not strictly necessary. However, in this paper we stick to the causal
 185 prior since it provides an explicit way of combining experimental and observational data.

186 For simplicity, in the sequel we use $\mathbf{b} = (\mathbf{v}, \mathbf{c})$ to denote in compact form the values of the
 187 variables in the base set $\mathbf{I} = \mathbf{v}$ and $\mathbf{C} = \mathbf{c}$. Using *do*-calculus we can compute $\hat{f}(\mathbf{b}) = \hat{f}(\mathbf{v}, \mathbf{c}) =$
 188 $\hat{\mathbb{E}}[Y | \text{do}(\mathbf{I} = \mathbf{v}), \mathbf{c}]$ and $\hat{\sigma}(\mathbf{b}) = \hat{\sigma}(\mathbf{v}, \mathbf{c}) = \hat{\mathbb{V}}[Y | \text{do}(\mathbf{I} = \mathbf{v}), \mathbf{c}]^{1/2}$ where $\hat{\mathbb{V}}$ and $\hat{\mathbb{E}}$ represent the
 189 variance and expectation of the causal effects estimated from \mathcal{D}^O . The *causal prior* $f(\mathbf{b}) \sim$
 190 $\mathcal{GP}(m(\mathbf{b}), K(\mathbf{b}, \mathbf{b}'))$ is defined to have prior mean and variance given by $m(\mathbf{b}) = \hat{f}(\mathbf{b})$ and
 191 $K(\mathbf{b}, \mathbf{b}') = k_{\text{RBF}}(\mathbf{b}, \mathbf{b}') + \hat{\sigma}(\mathbf{b})\hat{\sigma}(\mathbf{b}')$ respectively. The term $k_{\text{RBF}}(\mathbf{b}, \mathbf{b}') := \sigma_f^2 \exp(-\|\mathbf{b} - \mathbf{b}'\|^2 / 2l^2)$
 192 denotes the radial basis function (RBF) kernel, which is added to provide more flexibility to the model.

193 *Step 2, propagating the distribution to all elements in \mathbf{T} :* In Section 3.1 we showed how, $\forall \mathbf{X}_s \in$
 194 $\mathcal{P}(\mathbf{X})$, $t_s(\mathbf{x}) = L_s(f)(\mathbf{x})$ with f given by the intervention function defined in Theorem 3.1. By
 195 linearity of the causal operator, placing a GP prior on f induces a well-defined joint GP prior
 196 distribution on \mathbf{T} . In particular, for each $\mathbf{X}_s \in \mathcal{P}(\mathbf{X})$, we have $t_s(\mathbf{x}) \sim \mathcal{GP}(m_s(\mathbf{x}), k_s(\mathbf{x}, \mathbf{x}'))$ with:

$$m_s(\mathbf{x}) = \int \cdots \int m(\mathbf{b}) \pi_s(\mathbf{x}, \mathbf{b}_s) d\mathbf{b}_s \quad (4)$$

$$k_s(\mathbf{x}, \mathbf{x}') = \int \cdots \int K(\mathbf{b}, \mathbf{b}') \pi_s(\mathbf{x}, \mathbf{b}_s) \pi_s(\mathbf{x}', \mathbf{b}'_s) d\mathbf{b}_s d\mathbf{b}'_s. \quad (5)$$

197 where $\mathbf{b}_s = (\mathbf{v}_s^N, \mathbf{c})$ is the subset of \mathbf{b} including only the \mathbf{v} values corresponding to the set \mathbf{I}_s^N .

198 Let D be a finite set of inputs for the functions in \mathbf{T} , that is $D = \bigcup_s \{\mathbf{x}_{si}\}_{i=1}^M$. \mathbf{T} computed in D fol-
 199 lows a multivariate Gaussian distribution that is $\mathbf{T}^D \sim \mathcal{N}(m_{\mathbf{T}}(D), K_{\mathbf{T}}(D, D))$ with $K_{\mathbf{T}}(D, D) =$
 200 $(K_{\mathbf{T}}(\mathbf{x}, \mathbf{x}'))_{\mathbf{x} \in D, \mathbf{x}' \in D}$ and $m_{\mathbf{T}}(D) = (m_{\mathbf{T}}(\mathbf{x}))_{\mathbf{x} \in D}$. In particular, for two generic data points
 201 $\mathbf{x}_{si}, \mathbf{x}_{s'j} \in D$ with s and s' denoting two *distinct* functions we have $m_{\mathbf{T}}(\mathbf{x}_{si}) = \mathbb{E}[t_s(\mathbf{x}_i)] = m_s(\mathbf{x}_i)$
 202 and $K_{\mathbf{T}}(\mathbf{x}_{si}, \mathbf{x}_{s'j}) = \text{Cov}[t_s(\mathbf{x}_i), t_{s'}(\mathbf{x}_j)]$.

203 When computing the covariance function across intervention sets and intervention levels we differen-
 204 tiate between two cases. When both t_s and $t_{s'}$ are different from f , we have:

$$\text{Cov}[t_s(\mathbf{x}_i), t_{s'}(\mathbf{x}_j)] = \int \cdots \int K(\mathbf{b}, \mathbf{b}') \pi_s(\mathbf{x}_i, \mathbf{b}_s) \pi_{s'}(\mathbf{x}_j, \mathbf{b}'_{s'}) d\mathbf{b}_s d\mathbf{b}'_{s'}.$$

		Interventional data	
		No	Yes
Observational data	No		
	Yes		
		Single-task	Multi-task
		GP	DAG-GP
	Mechanistic model	$p(\mathbf{T}) = \prod_s p(t_s(\mathbf{x}))$ $t_s(\mathbf{x}) \sim \mathcal{GP}(0, K_{RBF}(\mathbf{x}, \mathbf{x}'))$	$p(\mathbf{T}) = \prod_s p(t_s(\mathbf{x}) f)$ $t_s(\mathbf{x}) = \int f(\mathbf{b})\pi_s(\mathbf{x}, \mathbf{b}_s)d\mathbf{b}_s$ $f(\mathbf{b}) \sim \mathcal{GP}(0, K_{RBF}(\mathbf{b}, \mathbf{b}'))$
	do-calculus	GP^+ $p(\mathbf{T}) = \prod_s p(t_s(\mathbf{x}))$ $t_s(\mathbf{x}) \sim \mathcal{GP}(m^+(\mathbf{x}), K^+(\mathbf{x}, \mathbf{x}'))$	DAG-GP^+ $p(\mathbf{T}) = \prod_s p(t_s(\mathbf{x}) f)$ $t_s(\mathbf{x}) = \int f(\mathbf{b})\pi_s(\mathbf{x}, \mathbf{b}_s)d\mathbf{b}_s$ $f(\mathbf{b}) \sim \mathcal{GP}(m^+(\mathbf{b}), K^+(\mathbf{b}, \mathbf{b}'))$

Figure 3: Models for learning the intervention functions \mathbf{T} defined on a DAG. The *do*-calculus allows estimating \mathbf{T} when only the observational data is available. When the interventional data is also available, one can use a single-task model (denoted by GP) or a multi-task model (denoted by DAG-GP). When both data types are available one can combine them using the causal prior parameters represented by $m^+(\cdot)$ and $k^+(\cdot, \cdot)$. The resulting models are denoted by GP^+ and DAG-GP^+ .

205 If one of the two functions equals f , this expression further reduces to:

$$\text{Cov}[t_s(\mathbf{x}_i), t_{s'}(\mathbf{x}_j)] = \int K(\mathbf{b}, \mathbf{b}')\pi_{s'}(\mathbf{x}_j, \mathbf{b}'_{s'})d\mathbf{b}'_{s'}.$$

206 Note that the integrating measures $\pi_s(\cdot, \cdot)$ and $\pi_{s'}(\cdot, \cdot)$ allow to compute the covariance between
 207 points that are defined on spaces on possibly different dimensionality, *a scenario that traditional*
 208 *multi-output GP models are unable to handle*. The prior $p(\mathbf{T})$ enables to merge different data types
 209 and to account for the natural correlation structure among interventions defined by the topology
 210 of the DAG. For this reason we call this formulation the DAG-GP model. The parameters in Eqs.
 211 (4)–(5) can be computed in closed form only when $K(\mathbf{b}, \mathbf{b}')$ is an RBF kernel and the integrating
 212 measures are assumed to be Gaussian distributions. In all other cases, one needs to resort to numerical
 213 approximations e.g. Monte Carlo integration in order to compute the parameters of each $t_s(\mathbf{x})$.

214 **Posterior distribution on \mathbf{T} :** The posterior distribution $p(\mathbf{T}^D|\mathcal{D}^I)$ can be derived analytically via
 215 standard GP updates. For any set D , $p(\mathbf{T}^D|\mathcal{D}^I)$ will be Gaussian with parameters $m_{\mathbf{T}|\mathcal{D}^I}(D) =$
 216 $m_{\mathbf{T}}(D) + K_{\mathbf{T}}(D, \mathbf{X}^I)[K_{\mathbf{T}}(\mathbf{X}^I, \mathbf{X}^I) + \sigma^2\mathbf{I}]^{-1}(\mathbf{T}^I - m_{\mathbf{T}}(\mathbf{X}^I))$ and $K_{\mathbf{T}|\mathcal{D}^I}(D, D) = K_{\mathbf{T}}(D, D) -$
 217 $K_{\mathbf{T}}(D, \mathbf{X}^I)[K_{\mathbf{T}}(\mathbf{X}^I, \mathbf{X}^I) + \sigma^2\mathbf{I}]^{-1}K_{\mathbf{T}}(\mathbf{X}^I, D)$. See Fig. 2 for an illustration of the DAG-GP model.

218 4 A helicopter view

219 Different variations of the DAG-GP model can be considered depending on the availability of both
 220 observational \mathcal{D}^O and interventional data \mathcal{D}^I (Fig. 3). Our goal here is not to be exhaustive, nor
 221 prescriptive, but to help to give some perspective. When \mathcal{D}^I is not available *do*-calculus is the only
 222 way to learn \mathbf{T} , which in turns requires \mathcal{D}^O . When both data types are not available, learning \mathbf{T} via a
 223 probabilistic model is not possible unless the causal effects can be transported from an alternative
 224 population. In this case mechanistic models based on physical knowledge of the process under
 225 investigation are the only option. When \mathcal{D}^I is available one can consider a single task or a multi-task
 226 model. If f does not exist, a single GP model needs to be considered for each intervention function.
 227 Depending on the availability of \mathcal{D}^O , integrating observational data into the prior distribution (denoted
 228 by GP^+) or adopting a standard prior (denoted by GP) are the two alternatives. In both cases, the
 229 experimental information is not shared across functions and learning \mathbf{T} requires intervening on all
 230 sets in $\mathcal{P}(\mathbf{X})$. When instead f exists, DAG-GP can be used to transfer interventional information and,
 231 depending on \mathcal{D}^O , also incorporating observational information a priori (DAG-GP^+).

232 5 Experiments

233 This section evaluates the performance of the DAG-GP model on two synthetic settings and on a real
 234 world healthcare application (Fig. 4). We first learn \mathbf{T} with fixed observational and interventional data
 235 (§5.1) and then use the DAG-GP model to solve active learning (AL) (§5.2) and Bayesian Optimization
 236 (BO) (§5.3)⁴. Implementation details are given in the supplement.

⁴Code and data for all the experiments will be provided.

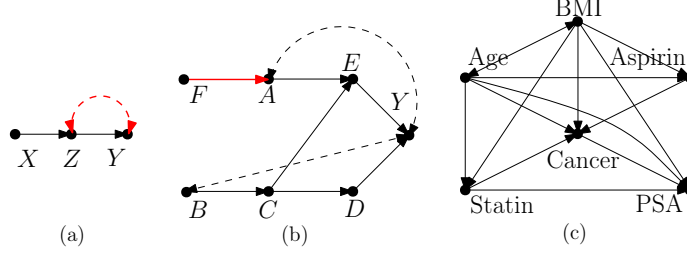


Figure 4: Examples of DAGs (in black) for which f exists and the DAG-GP model can be formulated. The red edges, if added, prevent the identification of f making the transfer via DAG-GP not possible.

Table 1: RMSE performances across 10 initializations of \mathcal{D}^I . See Fig. 3 for details on the compared methods. *do* stands for the *do*-calculus. N is the size of \mathcal{D}^O . Standard errors in brackets.

		$N = 30$					$N = 100$				
		DAG-GP ⁺	DAG-GP	GP ⁺	GP	<i>do</i>	DAG-GP ⁺	DAG-GP	GP ⁺	GP	<i>do</i>
DAG1		0.46 (0.06)	0.57 (0.09)	0.60 (0.2)	0.77 (0.27)	0.70 -	0.43 (0.05)	0.57 (0.08)	0.45 (0.05)	0.77 (0.27)	0.52 -
DAG2		0.44 (0.1)	0.45 (0.13)	0.62 (0.10)	1.26 (0.11)	1.40 -	0.36 (0.09)	0.41 (0.12)	0.58 (0.07)	1.28 (0.11)	1.41 -
DAG3		0.05 (0.04)	0.44 (0.12)	0.23 (0.03)	0.89 (0.23)	0.18 -	0.06 (0.04)	0.44 (0.12)	0.48 (0.06)	0.89 (0.23)	0.23 -

Baselines: We run our algorithm both with (DAG-GP⁺) and without (DAG-GP) causal prior and compare against the alternative models described in Fig. 3. Note that we do not compare against alternative multi-task GP models because, as mentioned in Section 1.2, the models existing in the literature cannot deal with functions defined on different inputs spaces and thus can not be straightforwardly applied to our problem.

Performance measures: We run all models with different initialisation of \mathcal{D}^I and different sizes of \mathcal{D}^O . We report the root mean square error (RMSE) performances together with standard errors across replicates. For the AL experiments we show the RMSE evolution as the size of \mathcal{D}^I increases. For the BO experiments we report the convergence performances to the global optimum.

5.1 Learning T from data

We test the algorithm on the DAGs in Fig. 4 and refer to them as (a) DAG1, (b) DAG2 and (c) DAG3. DAG3 is taken from [32] and [13] and is used to model the causal effect of statin drugs on the levels of prostate specific antigen (PSA). We consider the nodes $\{A, C\}$ in DAG2 and $\{\text{age, BMI, cancer}\}$ in DAG3 to be non-manipulative. We set the size of \mathcal{D}^I to $5 \times |\mathbf{T}|$ for DAG1 ($|\mathbf{T}| = 2$), to $3 \times |\mathbf{T}|$ for DAG2 ($|\mathbf{T}| = 6$) and to $|\mathbf{T}|$ for DAG3 ($|\mathbf{T}| = 3$). As expected, GP⁺ outperforms GP incorporating the information in \mathcal{D}^O (Tab. 1). Interestingly, GP⁺ also outperforms DAG-GP in DAG3 when $N = 30$ and in DAG1 when $N = 100$. This depends on the effect that \mathcal{D}^O has, through its size N and its coverage of the interventional domains, on both the causal prior and the estimation of the integrating measures. Lower N and coverage imply not only a less precise estimation of the *do*-calculus but also a worse estimation of the integrating measures and thus a lower transfer of information. Higher N and coverage imply more accurate estimation of the causal prior parameters and enhanced transfer of information across experiments. In addition, the way \mathcal{D}^O affects the performance results it's specific to the DAG structure and to the distribution of the exogenous variables in the SCM. More importantly, Tab. 1 shows how DAG-GP⁺ consistently outperforms all competing methods by successfully integrating different data sources and transferring interventional information across functions in \mathbf{T} . Differently from competing methods, these results holds across different N and \mathcal{D}^I values making DAG-GP⁺ a robust default choice for any application.

5.2 DAG-GP as surrogate model in Active Learning

The goal of AL is to design a sequence of function evaluations to perform in order to learn a target function as quickly as possible. We run DAG-GP within the AL algorithm proposed by [20] and select observations based on the Mutual Information (MI) criteria extended to a multi-task setting (see §5.2

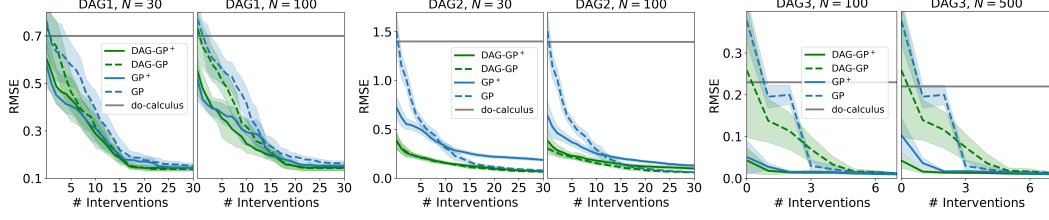


Figure 5: AL results. Convergence of the RMSE performance across functions in \mathbf{T} and across replicates as more experiments are collected. DAG-GP⁺ gives our algorithm with the causal prior while DAG-GP is our algorithm with a standard prior. # interventions is the number of experiments for each \mathbf{X}_s . Shaded areas give \pm standard deviation. See Fig. 3 for details on the compared methods.

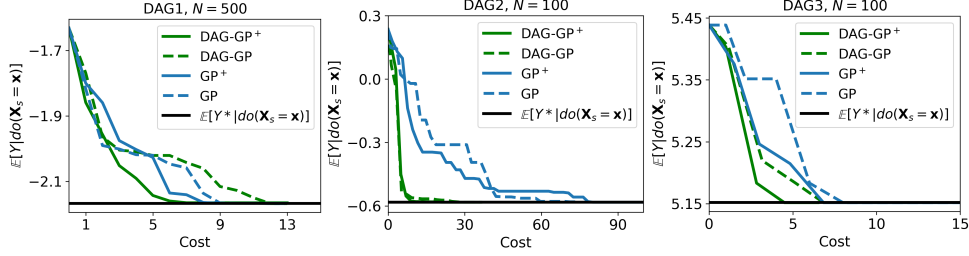


Figure 6: BO results. Convergence of the CBO algorithm to the global optimum ($\mathbb{E}[Y^* | \text{do}(\mathbf{X}_s = \mathbf{x})]$) when our algorithm is used as a surrogate model with (DAG-GP⁺) and without (DAG-GP) the causal prior. See the supplement for standard deviations across replicates.

in the supplement for details). Fig. 5 shows the RMSE performances as more interventional data are collected. Across different N settings, DAG-GP⁺ converges to the lowest RMSE performance faster than competing methods by collecting evaluations in areas where: (i) \mathcal{D}^O does not provide information and (ii) the predictive variance is not reduced by the experimental information transferred from the other interventions. As mentioned before, \mathcal{D}^O impacts on the causal prior parameters via the *do*-calculus computations. When the latter are less precise, because of lower N or lower coverage of the interventional domains, the model variances for DAG-GP⁺ or GP⁺ are inflated. Therefore, when DAG-GP⁺ or GP⁺ are used as surrogate models, the interventions are collected mainly in areas where \mathcal{D}^O is not observed thus slowing down the exploration of the interventional domains and the convergence to the minimum RMSE (Fig. 5 DAG2, $N = 100$).

5.3 DAG-GP as surrogate model in Bayesian optimization

The goal of BO is to optimize a function which is costly to evaluate and for which an explicit functional form is not available by making a series of function evaluations. We use DAG-GP within the CBO algorithm proposed by [1] (Fig. 6 right plot) where a modified version of the expected improvement is used as an acquisition functions to explore a set of intervention functions. We compare DAG-GP against the single-task models used in [1]. We found DAG-GP to significantly speed up the convergence of CBO to the global optimum both with and without the causal prior.

6 Conclusions

This paper addresses the problems of modelling the correlation structure of a set of intervention functions defined on the DAG of a causal model. We propose the DAG-GP model, which is based on a theoretical analysis of the DAG structure, and allows to share experimental information across interventions while integrating observational and interventional data via *do*-calculus. Our results demonstrate how DAG-GP outperforms competing approaches in term of fitting performances. In addition, our experiments show how integrating decision making algorithms with the DAG-GP model is crucial when designing optimal experiments as DAG-GP accounts for the uncertainty reduction obtained by transferring interventional data. Future work will extend the DAG-GP model to allow for transfer of experimental information *across* environments whose DAGs are partially different. In addition, we will focus on combining the proposed framework with a causal discovery algorithm so as to account for uncertainty in the graph structure.

7 Broader Impact

Computing causal effects is an integral part of scientific inquiry, spanning a wide range of questions such as understanding behaviour in online systems, assessing the effect of social policies, or investigating the risk factors for diseases. By combining the theory of causality with machine learning techniques, Causal Machine Learning algorithms have the potential to highly impact society and businesses by answering what-if questions, enabling policy-evaluation and allowing for data-driven decision making in real-world contexts. The algorithm proposed in this paper falls into this category and focuses on addressing causal questions in a fast and accurate way. As shown in the experiments, when used within decision making algorithms, the DAG-GP model has the potential to speed up the learning process and to enable optimal experimentation decisions by accounting for the multiple causal connections existing in the process under investigation and their cross-correlation. Our algorithm can be used by practitioners in several domains. For instance, it can be used to learn about the impact of environmental variables on coral calcification [12] or to analyse the effects of drugs on cancer antigens [13]. In terms of methodology, while the DAG-GP model represents a step towards a better model for automated decision making, it is based on the crucial assumption of knowing the causal graph. Learning the intervention functions of an incorrect causal graph might lead to incorrect inference and sub-optimal decisions. Therefore, more work needs to be done to account for the uncertainty in the graph structure.

References

- [1] Aglietti, V., Lu, X. L., Paleyes, A., and González, J. (2020). Causal Bayesian Optimization. In *Artificial Intelligence and Statistics*.
- [2] Alaa, A. M. and Van der Schaar, M. (2017). Bayesian inference of individualized treatment effects using multi-task Gaussian processes. In *Advances in Neural Information Processing Systems*, pages 3424–3432.
- [3] Álvarez, M., Luengo, D., and Lawrence, N. D. (2009). Latent force models. In *Artificial Intelligence and Statistics*, pages 9–16.
- [4] Álvarez, M. A., Rosasco, L., Lawrence, N. D., et al. (2012). Kernels for vector-valued functions: A review. *Foundations and Trends® in Machine Learning*, 4(3):195–266.
- [5] Bareinboim, E., Forney, A., and Pearl, J. (2015). Bandits with unobserved confounders: A causal approach. In *Advances in Neural Information Processing Systems*, pages 1342–1350.
- [6] Bareinboim, E. and Pearl, J. (2012). Causal inference by surrogate experiments: z-identifiability. *arXiv preprint arXiv:1210.4842*.
- [7] Bareinboim, E. and Pearl, J. (2013). Meta-transportability of causal effects: A formal approach. In *Proceedings of the 16th International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 135–143.
- [8] Bareinboim, E. and Pearl, J. (2014). Transportability from multiple environments with limited experiments: Completeness results. In *Advances in neural information processing systems*, pages 280–288.
- [9] Bonilla, E. V., Chai, K. M., and Williams, C. (2008). Multi-task Gaussian process prediction. In *Advances in neural information processing systems*, pages 153–160.
- [10] Buesing, L., Weber, T., Zwols, Y., Racaniere, S., Guez, A., Lespiau, J.-B., and Heess, N. (2018). Woulda, coulda, shoulda: Counterfactually-guided policy search. *arXiv preprint arXiv:1811.06272*.
- [11] Cochran, W. and Cox, G. (1957). Experimental design. John Wiley and Sons. Inc., New York, NY.
- [12] Courtney, T. A., Lebrato, M., Bates, N. R., Collins, A., De Putron, S. J., Garley, R., Johnson, R., Molinero, J.-C., Noyes, T. J., Sabine, C. L., et al. (2017). Environmental controls on modern scleractinian coral and reef-scale calcification. *Science advances*, 3(11):e1701356.

- [13] Ferro, A., Pina, F., Severo, M., Dias, P., Botelho, F., and Lunet, N. (2015). Use of statins and serum levels of prostate specific antigen. *Acta Urológica Portuguesa*, 32(2):71–77.
- [14] Foerster, J. N., Farquhar, G., Afouras, T., Nardelli, N., and Whiteson, S. (2018). Counterfactual multi-agent policy gradients. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- [15] Galles, D. and Pearl, J. (2013). Testing identifiability of causal effects. *arXiv preprint arXiv:1302.4948*.
- [16] Greenewald, K., Katz, D., Shanmugam, K., Magliacane, S., Kocaoglu, M., Adsera, E. B., and Bresler, G. (2019). Sample efficient active learning of causal trees. In *Advances in Neural Information Processing Systems*, pages 14279–14289.
- [17] Guo, R., Cheng, L., Li, J., Hahn, P. R., and Liu, H. (2018). A survey of learning causality with data: Problems and methods. *arXiv preprint arXiv:1809.09337*.
- [18] Hauser, A. and Bühlmann, P. (2014). Two optimal strategies for active learning of causal models from interventional data. *International Journal of Approximate Reasoning*, 55(4):926–939.
- [19] He, Y.-B. and Geng, Z. (2008). Active learning of causal networks with intervention experiments and optimal designs. *Journal of Machine Learning Research*, 9(Nov):2523–2547.
- [20] Krause, A., Singh, A., and Guestrin, C. (2008). Near-optimal sensor placements in Gaussian processes: Theory, efficient algorithms and empirical studies. *Journal of Machine Learning Research*, 9(Feb):235–284.
- [21] Lattimore, F., Lattimore, T., and Reid, M. D. (2016). Causal bandits: Learning good interventions via causal inference. In *Advances in Neural Information Processing Systems*, pages 1181–1189.
- [22] Lee, S. and Bareinboim, E. (2018). Structural causal bandits: where to intervene? In *Advances in Neural Information Processing Systems*, pages 2568–2578.
- [23] Lee, S. and Bareinboim, E. (2019). Structural causal bandits with non-manipulable variables. Technical report, Technical Report R-40, Purdue AI Lab, Department of Computer Science, Purdue.
- [24] Lu, C., Schölkopf, B., and Hernández-Lobato, J. M. (2018). Deconfounding reinforcement learning in observational settings. *arXiv preprint arXiv:1812.10576*.
- [25] Magliacane, S., van Ommen, T., Claassen, T., Bongers, S., Versteeg, P., and Mooij, J. M. (2018). Domain adaptation by using causal inference to predict invariant conditional distributions. In *Advances in Neural Information Processing Systems*, pages 10846–10856.
- [26] Pearl, J. (1995). Causal diagrams for empirical research. *Biometrika*, 82(4):669–688.
- [27] Pearl, J. (2000). *Causality: models, reasoning and inference*, volume 29. Springer.
- [28] Pearl, J. and Bareinboim, E. (2011). Transportability of causal and statistical relations: A formal approach. In *Twenty-fifth AAAI conference on artificial intelligence*.
- [29] Rasmussen, C. E. (2003). Gaussian processes in machine learning. In *Summer School on Machine Learning*, pages 63–71. Springer.
- [30] Rojas-Carulla, M., Schölkopf, B., Turner, R., and Peters, J. (2018). Invariant models for causal transfer learning. *The Journal of Machine Learning Research*, 19(1):1309–1342.
- [31] Rubenstein, P. K., Tolstikhin, I., Hennig, P., and Schölkopf, B. (2017). Probabilistic active learning of functions in structural causal models. *arXiv preprint arXiv:1706.10234*.
- [32] Thompson, C. (2019). Causal graph analysis with the causalgraph procedure. <https://www.sas.com/content/dam/SAS/support/en/sas-global-forum-proceedings/2019/2998-2019.pdf>.

- 389 [33] Ye, C., Butler, L., Bartek, C., Iangurazov, M., Lu, Q., Gregory, A., Girolami, M., and Middleton,
390 C. (2019). A digital twin of bridges for structural health monitoring. In *12th International*
391 *Workshop on Structural Health Monitoring 2019*. Stanford University.
- 392 [34] Zhang, K., Schölkopf, B., Muandet, K., and Wang, Z. (2013). Domain adaptation under target
393 and conditional shift. In *International Conference on Machine Learning*, pages 819–827.