

ROBUST AGENTS LEARN CAUSAL WORLD MODELS

Jonathan Richens*
Google DeepMind

Tom Everitt
Google DeepMind

ABSTRACT

It has long been hypothesised that causal reasoning plays a fundamental role in robust and general intelligence. However, it is not known if agents must learn causal models in order to generalise to new domains, or if other inductive biases are sufficient. We answer this question, showing that any agent capable of satisfying a regret bound for a large set of distributional shifts must have learned an approximate causal model of the data generating process, which converges to the true causal model for optimal agents. We discuss the implications of this result for several research areas including transfer learning and causal inference.

1 INTRODUCTION

What capabilities are necessary for general intelligence (Legg & Hutter, 2007)? One candidate is causal reasoning, which plays a foundational role in human cognition (Gopnik et al., 2007; Sloman & Lagnado, 2015). It has even been argued that human-level AI is impossible without causal reasoning (Pearl, 2018). However, recent years have seen the development of agents that do not explicitly learn or reason on causal models, but nonetheless are capable of adapting to a wide range of environments and tasks (Reed et al., 2022; Team et al., 2023; Brown et al., 2020).

This raises the question, do agents have to learn causal models in order to adapt to new domains, or are other inductive biases sufficient? To answer this question, we have to be careful not to assume that agents use causal assumptions a priori. For example, transportability theory determines what causal knowledge is necessary for transfer learning when all assumptions on the data generating process (inductive biases) can be expressed as constraints on causal structure (Bareinboim & Pearl, 2016). However, deep learning algorithms can exploit a much larger set of inductive biases (Neyshabur et al., 2014; Battaglia et al., 2018; Rahaman et al., 2019; Goyal & Bengio, 2022) which in many real-world tasks may be sufficient to identify low regret policies without requiring causal knowledge.

The main result of this paper is to answer this question by showing that,

Any agent capable of adapting to a sufficiently large set of distributional shifts must have learned a causal model of the data generating process.

Here, adapting to a distributional shift means learning a policy that satisfies a regret bound following an intervention on the data generating process—for example, changing the distribution of features or latent variables. It is known that a causal model of the data generating process can be used to identify regret-bounded policies following a distributional shift (sufficiency), with more accurate models allowing lower regret policies to be found. We prove the converse (necessity)—given regret-bounded policies for a large set of distributional shifts, we can learn an approximate causal model of the data generating process, with the approximation becoming exact for optimal policies. Hence, learning a causal model of the data generating process is both necessary and sufficient for robust adaptation.

This equivalence has consequences for a number of fields and questions. For one, it implies that causal identification laws also constrain domain adaptation. For example, we show that adapting to covariate and label shifts is only possible if the causal relations between features and labels can be identified from the training data—a non-trivial causal discovery problem. This provides further theoretical justification for causal representation learning (Schölkopf et al., 2021), showing that learning causal representations is necessary for achieving strong robustness guarantees. Our result also implies that we can learn causal models from adaptive agents. We demonstrate this by solving a causal discovery

*jonrichens@deepmind.com

task on synthetic data by observing the policy of a regret-bounded agent under distributional shifts. More speculatively, our results suggest that causal models could play a role in emergent capabilities. Agents trained to minimise a loss function across many domains are incentivized to learn a causal world model, which could enable them to solve a much larger set of decision tasks they were not explicitly trained on.

Outline of paper. In Section 2 we introduce concepts from causality and decision theory used to derive our results. We present our main theoretical results in Section 3 and discuss their interpretation in terms of adaptive agents, transfer learning and causal inference. In Section 4 we discuss limitations, as well as implications for a number of fields and open questions. In section 5 we discuss related work including transportability (Bareinboim & Pearl, 2016) and the causal hierarchy theorem (Bareinboim et al., 2022), and recent empirical work on emergent world models. In Appendix B we describe experiments applying our theoretical results to causal discovery problems.

2 PRELIMINARIES

2.1 CAUSAL MODELS

We use capital letters for random variables V , and lower case for their values $v \in \text{dom}(V)$. For simplicity, we assume each variable has a finite number of possible values, $|\text{dom}(V)| < \infty$. Bold face denotes sets of variables $\mathbf{V} = \{V_1, \dots, V_n\}$, and their values $\mathbf{v} \in \text{dom}(\mathbf{V}) = \times_i \text{dom}(V_i)$. A probabilistic model specifies the joint distribution $P(\mathbf{V})$ over a set of variables \mathbf{V} . These models can support associative queries, for example $P(\mathbf{Y} = \mathbf{y} \mid \mathbf{X} = \mathbf{x})$ for $\mathbf{X}, \mathbf{Y} \subseteq \mathbf{V}$. Interventions describe external changes to the data generating process (and hence changing the joint distribution), for example a *hard* intervention $\text{do}(\mathbf{X} = \mathbf{x})$ describes forcing the set of variables $\mathbf{X} \subseteq \mathbf{V}$ to take value \mathbf{x} . This generates a new distribution $P(\mathbf{V} \mid \text{do}(\mathbf{X} = \mathbf{x})) = P(\mathbf{V}_x)$ where \mathbf{V}_x refers to the variables \mathbf{V} following this intervention. The power of causal models is that they specify not only $P(\mathbf{V})$ but also the distribution of \mathbf{V} under all interventions, and hence these models can be used to evaluate both associative and interventional queries e.g. $P(\mathbf{Y} = \mathbf{y} \mid \text{do}(\mathbf{X} = \mathbf{x}))$.

For the derivation of our results we focus on a specific class of causal models—causal Bayesian networks (CBNs). There are several alternative models and formalisms that are studied in the literature, including structural equation models (Pearl, 2009) and the Neyman-Rubin causal models (Rubin, 2005), and results can be straightforwardly adapted to these.

Definition 1 (Bayesian networks). A Bayesian network $M = (G, P)$ over a set of variables $\mathbf{V} = \{V_1, \dots, V_n\}$ is a joint probability distribution $P(\mathbf{V})$ that factors according to a directed acyclic graph (DAG) G , i.e. $P(V_1, \dots, V_n) = \prod_{i=1}^n P(V_i \mid \mathbf{Pa}_{V_i})$, where \mathbf{Pa}_{V_i} are the parents of V_i in G .

A Bayesian network is *causal* if the graph G captures the causal relationships between the variables or, formally, if the result of any intervention $\text{do}(\mathbf{X} = \mathbf{x})$ for $\mathbf{X} \subseteq \mathbf{V}$ can be computed from the truncated factorisation formula:

$$P(\mathbf{v} \mid \text{do}(\mathbf{x})) = \begin{cases} \prod_{i: v_i \notin \mathbf{x}} P(v_i \mid \mathbf{pa}_{v_i}) & \text{if } \mathbf{v} \text{ consistent with } \mathbf{x} \\ 0 & \text{otherwise.} \end{cases}$$

More generally, a *soft* intervention $\sigma_{v_i} = P'(V_i \mid \mathbf{Pa}_i^*)$ replaces the conditional probability distribution for V_i with a new distribution $P'(V_i \mid \mathbf{Pa}_i^*)$, possibly resulting in a new parent set $\mathbf{Pa}_i^* \neq \mathbf{Pa}_i$ as long as no cycles are introduced in the graph. We refer to σ_{v_i} as a *domain indicator* (Correa & Bareinboim, 2020) (it has also been called an environment index, Arjovsky et al., 2019). The updated distribution is denoted $P(\mathbf{v}; \sigma_{\mathbf{v}'}) = \prod_{i: v_i \in \mathbf{v}'} P'(v_i \mid \mathbf{pa}_{v_i}^*) \prod_{i: v_i \notin \mathbf{v}'} P(v_i \mid \mathbf{pa}_{v_i})$.

In general, soft interventions cannot be defined without knowledge of G . For example, the soft intervention $\sigma_Y = P'(y \mid x)$ is incompatible with the causal structure $Y \rightarrow X$ as it would induce a causal cycle. As our results are concerned with learning causal models (and hence causal structure), we focus our theoretical analysis on a subset of the soft interventions, *local interventions*, that are compatible with all causal structures and so can be used without tacitly assuming knowledge of G .

Definition 2 (Local interventions). *Local intervention σ on $V_i \in \mathbf{V}$ involves applying a map to the states of V_i that is not conditional on any other endogenous variables, $v_i \mapsto f(v_i)$. We use the notation $\sigma = \text{do}(V_i = f(v_i))$ (variable V_i is assigned the state $f(v_i)$). Formally, this is a soft intervention on V_i that transforms the conditional probability distribution as,*

$$P(v_i | \mathbf{pa}_i; \sigma) = \sum_{v'_i: f(v'_i)=v_i} P(v'_i | \mathbf{pa}_i) \quad (1)$$

Example: Hard interventions $\text{do}(V_i = v'_i)$ are local interventions where $f(v_i)$ is a constant function.

Example: Translations are local interventions as $\text{do}(V_i = v_i + k) = \text{do}(V_i = f(v_i))$ where $f(v_i) = v_i + k$. Examples include changing the position of objects in RL environments (Shah et al., 2022) and images (Engstrom et al., 2019).

Example: Logical NOT operation $X \mapsto \neg X$ for Boolean X is a local intervention.

We also consider stochastic interventions, noting that mixtures of local interventions can also be defined without knowledge of G . For example, adding noise to a variable $X = X + \epsilon$, $\epsilon \sim \mathcal{N}(0, 1)$, is a soft intervention on X described by a mixture over local interventions (translations).

Definition 3 (Mixtures of interventions). *A mixed intervention $\sigma^* = \sum_i p_i \sigma_i$ for $\sum p_i = 1$ performs intervention σ_i with probability p_i . Formally, $P(\mathbf{v} | \sigma^*) = \sum_i p_i P(\mathbf{v} | \sigma_i)$.*

2.2 DECISION TASKS

Decision tasks involve a decision maker (agent) choosing a policy so as to optimise an objective function (utility). To give a causal description of decision tasks we use the causal influence diagram (CID) formalism (Howard & Matheson, 2005; Everitt et al., 2021), which extend a CBN of the environment (chance) variables by introducing decision and utility nodes (see Figure 1 for examples). For simplicity we focus on tasks involving a single decision and a single utility function.

Definition 4 (Causal influence diagram). *A (single-decision, single-utility) causal influence diagram (CID) is a CBN $M = (G, P)$ where the variables \mathbf{V} are partitioned into decision, utility, and chance variables, $\mathbf{V} = (\{D\}, \{U\}, \mathbf{C})$. The utility variable is a real-valued function of its parents, $U(\mathbf{pa}_U)$.*

Single-decision single-utility CIDs can represent most decision tasks such as classification and regression as they specify what decision should be made ($d \in D$), based on what information (\mathbf{pa}_D), with objective ($\mathbb{E}[U]$). They can also describe some multi-decision tasks such as Markov decision processes¹. The utility is any real-valued function including standard loss and reward functions.

We assume that the environment is described by a set of random variables \mathbf{C} that interact via causal mechanisms², and where \mathbf{C} satisfies causal sufficiency (Pearl, 2009) (includes all common causes), noting that such a choice of \mathbf{C} always exists. We refer to the CBN over \mathbf{C} as the ‘true’ or ‘underlying’ CBN. Note we do not assume the agent has any knowledge of the underlying CBN, nor do we assume which variables in \mathbf{C} are observed or unobserved by the agent, beyond that the agent can observe $\mathbf{Pa}_D \subseteq \mathbf{C}$. We also assume knowledge of the utility function $U(\mathbf{Pa}_U)$.

The conditional probability distribution for the decision node $\pi(d | \mathbf{pa}_D)$ (the policy) is not a fixed parameter of the model but is set by the agent so as to maximise its expected utility, which for a policy π is $\mathbb{E}^\pi[U] = \mathbb{E}[U | \text{do}(D = \pi(\mathbf{pa}_D))]$. A policy π^* is *optimal* if it maximises $\mathbb{E}^{\pi^*}[U]$. Typically, agents do not behave optimally and incur some *regret* δ , which is the decrease in expected utility compared to an optimal policy $\delta := \mathbb{E}^{\pi^*}[U] - \mathbb{E}^\pi[U]$.

To simplify our theoretical analysis, we focus on a widely studied class of decision tasks where the agents decision does not causally influence the environment (e.g. Figure 1).

Assumption 1 (Unmediated decision task). $\text{Desc}_D \cap \text{Anc}_U = \emptyset$.

In unmediated decision tasks, the agent is provided some (partial) observations of the environment and chooses a policy, which is then evaluated using the utility function which is a function of the

¹Note Markov decision processes can be formulated as a single-decision single-utility CID, by modelling the choice of policy as a single decision and the cumulative discounted reward as a single utility variable.

²This assumption follows from Reichenbach (1956), and we discuss further in Appendix A.3

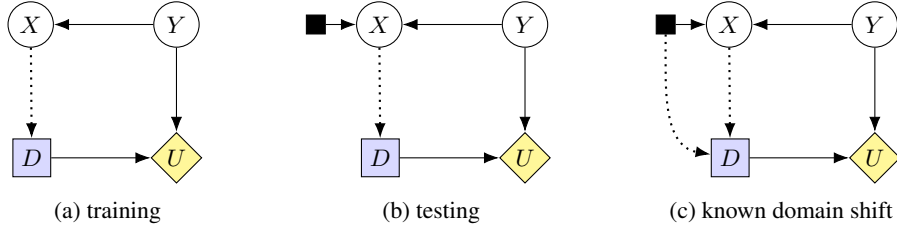


Figure 1: CID for a supervised learning task during (a) training and (b) testing following a distributional (covariate) shift (unsupervised domain adaptation, [Wilson & Cook, 2020](#)). The agent chooses a label prediction $D = \hat{Y}$ given features X , with the goal of minimising loss $U = -\text{Loss}(Y, \hat{Y})$. Decision variables are depicted as square nodes, chance variables as circular nodes and utilities as diamond nodes. Information edges (dashed) show the variables the agent conditions their policy on. In this example the labels cause the features $Y \rightarrow X$ (for examples where features cause labels see [Castro et al., 2020](#); [Schölkopf et al., 2012](#)). The black square (‘regime node’ ([Correa & Bareinboim, 2020](#))) in (b) and (c) denotes a distributional shift induced by an intervention on X . Diagram (c) depicts the idealised case where the agent knows what domain shift has occurred. By theorem [1](#), if the agent can return an optimal decision boundary for known covariate and label shifts, then it must have learned the CBN over $\mathcal{C} = \{X, Y\}$. Note that even if the agent has sufficient training data to learn $P(X, Y)$, the causal structure $Y \rightarrow X$ is in general non-identifiable given $P(X, Y)$ and so domain adaptation requires that the agent solves a non-trivial causal discovery problem.

environment state and the agent’s decision. Examples of unmediated decision tasks include prediction tasks such as classification and regression, whereas examples of *mediated* decision tasks that are not covered by our theorems include Markov decision processes where the agent’s decision (action) influences the utility via the environment state.

2.3 DISTRIBUTIONAL SHIFTS

We focus on generalisation that goes beyond the *iid* assumption, where agents are evaluated in domains that are distributionally shifted from the training environment. Distributional shifts can be changes to the environment (*domain shifts*), as in domain adaptation and domain generalisation ([Farahani et al., 2021](#); [Wilson & Cook, 2020](#)), or changes to the objective (*task shifts*) as in zero shot learning ([Xian et al., 2018](#)), in-context learning ([Brown et al., 2020](#)) and multi-task reinforcement learning ([Reed et al., 2022](#)). Our analysis focuses on domain shifts that involve changes to the causal data generating process, and hence can be modelled as interventions ([Schölkopf et al., 2021](#)). This does not assume that all shifts an agent will encounter can be modelled as interventions, but requires that the agent is *at least* capable of adapting to these shifts.

Examples of interventionally generated shifts include translating objects in images ([Engstrom et al., 2019](#)), noising inputs and adversarial robustness ([Hendrycks & Dietterich, 2019](#)), and changes to the initial conditions or transition function in Markov decision processes ([Peng et al., 2018](#)). Examples of shifts that are not naturally represented as interventions include changing the set of environment variables \mathcal{C} , and introducing selection biases ([Shen et al., 2018](#)). See Appendix [A.3](#) for discussion.

Our main results restrict to local *domain shifts*, which correspond to local interventions on the chance variables \mathcal{C} . We do not consider shifts that change the agent’s decision D , although we include shifts that drop inputs to the policy $\mathbf{Pa}_D \rightarrow \mathbf{Pa}'_D \subseteq \mathbf{Pa}_D$ (e.g. masking) as local interventions. We do not consider task shifts i.e. changing the utility function.

As we are interested in determining the capabilities necessary for domain adaptation, we restrict our attention to decision tasks where domain adaptation is non-trivial, i.e. where the optimal policy depends on the environment distribution $P(\mathcal{C} = c)$.

Assumption 2 (Domain dependence). *There exists $P(\mathcal{C} = c)$ and $P'(\mathcal{C} = c)$ compatible with M such that $\pi^* = \arg \max_{\pi} \mathbb{E}_{\mathcal{P}}^{\pi}[U]$ implies $\pi^* \neq \arg \max_{\pi} \mathbb{E}_{\mathcal{P}'}^{\pi}[U]$.*

Assumption 2 implies the existence of domain shifts that change the optimal policy. As shown in Appendix [A.1](#) domain independence holds if and only if $\mathbf{Anc}_U \subseteq \mathbf{Pa}_D$, i.e. when there are no latent variables (Lemma [1](#)). It is simple to show that in this case the optimal policy is invariant under all

domain shifts. Hence, Assumption 2 is equivalent to assuming the presence of latent variables that are strategically relevant (Koller & Milch, 2003).

3 EQUIVALENCE OF LEARNING POLICIES AND CAUSAL MODELS

We now present our results, showing that learning the underlying CBN is necessary and sufficient for learning regret bounded policies under a large set of domain shifts. In Section 3.2 we interpret these results by considering their consequences for transfer learning, causal inference and adaptive agents.

First we focus on the idealised case where we assume optimality. We show for almost all decision tasks the underlying CBN can be reconstructed given optimal policies for a large set of domain shifts.

Theorem 1. *For almost all CIDs $M = (G, P)$ satisfying Assumptions 1 and 2, we can identify the directed acyclic graph G and joint distribution P over all ancestors of the utility \mathbf{Anc}_U given $\{\pi_\sigma^*(d \mid \mathbf{pa}_D)\}_{\sigma \in \Sigma}$ where $\pi_\sigma^*(d \mid \mathbf{pa}_D)$ is an optimal policy in the domain σ and Σ is the set of all mixtures of local interventions. Proof in Appendix C*

The parameters $P(v_i \mid \mathbf{pa}_i)$, $U(\mathbf{pa}_U)$ of the underlying CBN define a parameter space and the condition for almost all CIDs means that the subset of the parameter space for which the Theorem 1 does not hold is Lebesgue measure zero (see Appendix A.2 for discussion). This condition is necessary because there exist finely-tuned environments for which the CBN cannot be identified given the agent’s policy due to variables $X \in \mathbf{Anc}_U$ that do not affect the expected utility. For example consider $X \rightarrow Y \rightarrow U$, $Y = \mathcal{N}(0, x)$ and $U = D + Y$, then changing X can only change the variance of U while leaving its expected value (and hence the optimal policy) constant. However, this only occurs for very specific choices of the parameters P and U .

In Appendix B we give a simplified overview of the proof with a worked example. We assume access to an oracle for optimal policies π_σ^* for any given local intervention on \mathcal{C} . Note, this assumes the agent is robust to distributional shifts on a causally sufficient set of variables \mathcal{C} , not that the set of variables the agent observes is causally sufficient. We devise an algorithm that queries this oracle with different mixtures of local interventions and identifies the mixtures for which the optimal policies changes. We then show that these critical mixtures identify the parameters of the CBN, specifying both the graph $G(\mathbf{Anc}_U)$ and the joint distribution $P(\mathbf{Anc}_U)$.

3.1 RELAXING THE ASSUMPTION OF OPTIMALITY

We now relax the assumption of optimality, considering the case where the policies π_σ satisfy a regret bound $\mathbb{E}^{\pi_\sigma}[U] \geq \mathbb{E}^{\pi_\sigma^*}[U] - \delta$. We show that for $\delta > 0$ we can recover an approximation of the environment CBN, with error that grows linearly in δ for $\delta \ll \mathbb{E}^{\pi^*}[U]$.

Theorem 2. *For almost all CIDs $M = (G, P)$ satisfying Assumptions 1 and 2 we can identify an approximate causal model $M' = (P', G')$ given $\{\pi_\sigma(d \mid \mathbf{pa}_D)\}_{\sigma \in \Sigma}$ where $\mathbb{E}^{\pi_\sigma}[U] \geq \mathbb{E}^{\pi_\sigma^*}[U] - \delta$ and Σ is the set of mixtures of local interventions. The parameters of M' satisfy $|P'(v_i \mid \mathbf{pa}_i) - P(v_i \mid \mathbf{pa}_i)| \leq \gamma(\delta) \forall V_i \in \mathbf{V}$ where $\gamma(0) = 0$ and $\gamma(\delta)$ grows linearly in δ for small regret $\delta \ll \mathbb{E}^{\pi^*}[U]$. Proof in Appendix D*

The worst-case error bounds $\gamma(\delta)$ for the parameter errors are detailed in Appendix D. For $\delta > 0$ it may not be possible to identify G perfectly as some weak causal relations cannot be resolved due to these error bounds. We describe in Appendix D how we can learn a sub-graph $G' \subseteq G$ that may exclude directed edges corresponding to weak causal relations.

Theorem 2 shows that we can learn a (sparse) approximate causal models of the data generating process from regret bounded policies under domain shifts, where the approximation becoming exact as $\delta \rightarrow 0$. In Appendix E we demonstrate learning the underlying CBN from regret-bounded policies using simulated data for randomly generated CIDs similar to Figure 1 and explore how the accuracy of the approximate CBN scales with the regret bound (Figure 3).

Finally, we prove sufficiency, i.e. that having an (approximate) causal model of the data generating process is sufficient to identify regret-bounded policies. The result is well-known for the non-approximate case (Bareinboim & Pearl, 2016).

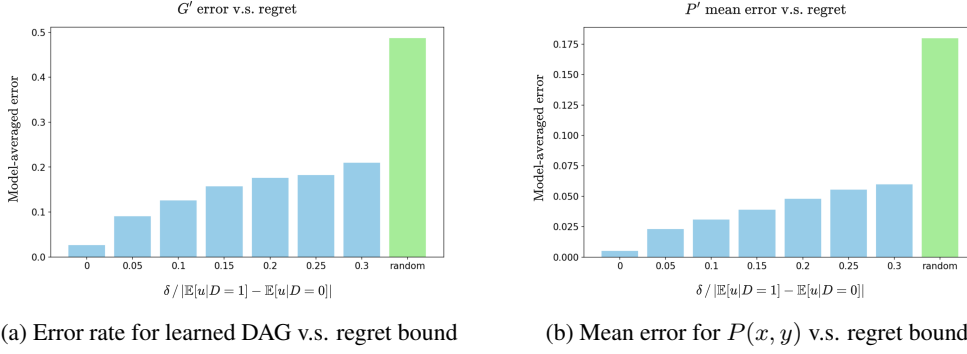


Figure 3: Comparing the model-average error rates for a) the learned DAG G' and b) learned joint distribution $P'(x, y)$, v.s. the (normalised) regret bound $\delta / |\mathbb{E}[u | D = 1] - \mathbb{E}[u | D = 0]|$. Average error taken over 1000 randomly generated environments with binary decision D and two binary latent variables X, Y . Comparison to error rate for random guess (green) See Appendix F for details.

Theorem 3. *Given the CBN $M = (P, G)$ that is causally sufficient we can identify optimal policies $\pi_\sigma^*(d | \mathbf{pa}_D)$ for any given U where $\mathbf{pa}_U \subseteq \mathcal{C}$ and for all soft interventions σ . Given an approximate causal model $M' = (P', G')$ for which $|P'(v_i | \mathbf{pa}_i) - P(v_i | \mathbf{pa}_i)| \leq \epsilon \ll 1$, we can identify regret-bounded policies where the regret δ grows linearly in ϵ . Proof in Appendix E.*

Together, Theorems 2 and 3 imply that learning an approximate causal model of the data generating process is necessary and sufficient for learning regret-bounded policies under local interventions.

3.2 INTERPRETATION

We interpret Theorems 1 to 3 through three lenses; transfer learning, adaptive agents and causal inference.

Transfer learning. In transfer learning (Zhuang et al., 2020), models are trained on a set of source domains and evaluated on held-out target domains where i) the data distribution differs from the source domains, and ii) the data available for training is restricted compared to the source domains (Wang et al., 2022). For example, in unsupervised domain adaptation the learner is restricted to samples of the input features from the target domains $\mathbf{pa}_D \sim P(\mathbf{pa}_D; \sigma)$, whereas in domain generalisation typically no data from the target domain is available during training (Farahani et al., 2021).

Let \mathcal{D}_S denote the training data from the source domains and \mathcal{D}_σ denote the training data available from a given target domain σ . Let there exist a transfer learning algorithm that returns a policy π_σ satisfying a regret bound for a given target domain σ , provided this training data. As π_σ is a function of the training data, then by Theorems 1 and 2 the existence of this algorithm implies that we can identify the underlying CBN from $\mathcal{D}_S \cup \{\mathcal{D}_\sigma\}_{\sigma \in \Sigma}$. To see that this imparts non-trivial constraints on the existence of the transfer learning algorithm, we can consider the following simple example.

Example: Consider the CID for the supervised learning task depicted in Figure 1. Let $\mathcal{D}_S = \{(x^i, y^i) \sim P(X, Y)\}_{i=1}^n$, so for sufficiently large n the agent can learn the $P(X, Y)$ from \mathcal{D}_S . However, $Y \rightarrow X$ must also be identifiable from the training data $\mathcal{D}_S \cup \{\mathcal{D}_\sigma\}_{\sigma \in \Sigma}$. In other words, the transfer learning problem contains a hidden causal discovery problem. If $\mathcal{D}_\sigma = \emptyset$ then $Y \rightarrow X$ must be identifiable from $P(x, y)$ alone, which is impossible unless the causal data generating process obeys additional assumptions (see for example Hoyer et al., 2008). If unlabelled features from the target domain are included in the training data $\mathcal{D}_\sigma = \{x^i \sim P(X; \sigma)\}_{i=1}^{n_\sigma}$, $Y \rightarrow X$ can in principle be identified as $P(X; \sigma_Y) \neq P(X)$.

Adaptive agents are goal-directed systems whose outputs are ‘moved by reasons’ (Dennett, 1989), meaning they choose an action because they expect it to achieve some desired outcome, and would act differently if they knew the consequences of their actions would be different. For example, a firm sets prices to maximise profit, and adapts pricing to changes in demand (Kenton et al., 2023).

Consider the transfer learning setting where an agent has to generalise to a target domain using only its previous experience (i.e. zero-shot adaptation), enabled by the fact that it has perfect knowledge of what domain shift has occurred $D_\sigma = \{\sigma\}$ (e.g. the agent conditions their policy on the domain indicator³ $\pi_\sigma = \pi(d \mid \mathbf{pa}_D, \sigma)$). If the policy π_σ satisfies a tight regret bound then by Theorem 2 we can reconstruct the underlying CBN from the agent’s policy alone (following the procedure described in Appendix C). Hence, any agent that is capable of adapting to known domain shifts has also learned a causal model of their environment.

Example: Doctors are one such agent, as they are expected to make low regret decisions under a wide range of known distributional shifts without re-training in the shifted environment. For example, consider a doctor tasked with risk-stratifying patients based on their signs and medical history. The doctor may be transferred to a new ward where patients have received a treatment (known distributional shift) that has a stochastic effect on latent variables (mixed intervention) such as curing diseases and causing side effects. The doctor cannot re-train in this new domain, e.g. taking random decisions and observing outcomes. To be capable of this, Theorem 2 implies the doctor must have learned the causal relations between the relevant latent variables—how the treatment affects diseases, how these diseases and their symptoms are causally related, and so on. Likewise, any medical AI that hopes to replicate this capability must have learned a similarly accurate causal model, and the better the agent’s performance the more accurate its causal model must be.

Where is the causal model? When we say this agent has learned the underlying CBN, we mean that the agent has learned the policy $\pi(d \mid \mathbf{pa}_D, \sigma)$ which is functionally equivalent to learning the underlying CBN, as any query that can be answered using the CBN can also be answered using $\pi(d \mid \mathbf{pa}_D, \sigma)$, which follows from the fact that the CBN is identified by $\pi(d \mid \mathbf{pa}_D, \sigma)$ (Theorem 1).

Causal inference. Theorem 1 can also be interpreted purely in terms of causal inference. We can compare to the causal hierarchy theorem (CHT) (Bareinboim et al. 2022), which states that an oracle for L1 queries (observational) is almost always insufficient to evaluate all L2 queries (interventional). Our Theorem 1 can be stated in an analogous way; an oracle for optimal policies under mixtures of local interventions $\Pi_\Sigma^* : \sigma \mapsto \pi^*(\sigma)$, can evaluate all L2 queries, which follows from the fact that the oracle identifies the underlying CBN which in turn identifies all L2 queries. Note Π_Σ^* is a strict subset of L2, and we describe a subset of L2 as being *L2-complete* if evaluating these queries is sufficient to evaluate all L2 queries. Hence Theorem 1 can be summarised as Π_Σ^* is L2-complete. It would be interesting in future work to determine what other strict subsets of L2 are L2-complete, as identifying these queries is sufficient to identify all interventional queries.

Why is this surprising? Firstly, we may expect the optimal policies to encode a relatively small number of causal relations, as they can be computed from $\mathbb{E}[u \mid d, \mathbf{pa}_D; \sigma]$, which describes the response of a single variable U to intervention σ . However, Theorem 1 shows that the optimal policies encode all causal and associative relations in \mathbf{Anc}_U , including causal relations between latent variables, for example $P(Y_x)$ for any $X, Y \subseteq \mathbf{Anc}_U$. Secondly, Theorems 2 and 3 combined imply that learning to generalise under domain shifts is *equivalent* to learning a causal model of the data generating process—problems that on the surface are conceptually distinct.

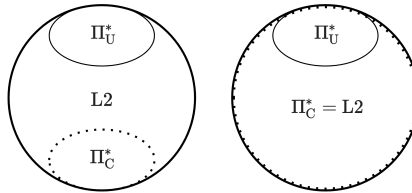


Figure 4: The set L2 (Bareinboim et al. 2022) contains all casual queries including the set of optimal policies under domain π_C^* and task shifts π_U^* . We show that π_C^* also contains L2, i.e. learning optimal policies under all shifts for a single utility U is sufficient to identify L2

³Domain indicator σ is equivalent to an environment index (Gupta et al. 2023; Arjovsky et al. 2019)

4 DISCUSSION

Here we discuss the consequences for several fields and open questions, as well as limitations.

Causal representation learning. Causal representation learning (CRL) aims to learn representations of data that capture unknown causal structure (Schölkopf et al., 2021), with the aim of exploiting causal invariances to achieve better generalisation across domains. Theorems 1 and 2 show that any method that enables generalisation across many domains necessarily involves learning an (approximate) causal model of the data generating process—i.e. a causal representation. Hence, our results provide theoretical justification for CRL by showing it is necessary for strong robustness guarantees.

Causal bounds on transfer learning. As described in Section 3.2, Theorems 1 and 2 imply fundamental causal constraints on certain transfer learning tasks. For example in the supervised learning task depicted in Figure 1, identifying regret-bounded policies under covariate and label shifts requires learning the causal relations between features and labels. Causal discovery problems such as this are well understood in many settings (Vowels et al., 2022), and in general identifying this causal structure (e.g. that $Y \rightarrow X$ in Figure 5a) is impossible without interventional data and/or additional assumptions. This connection allows us to convert (im)possibility results for causal discovery to (im)possibility results for transfer learning. Future work could explore this for smaller sets of distributional shifts and derive more general causal bounds on transfer learning.

Good regulator theorem. The good regulator theorem is often interpreted as saying that any good controller of a system must have a model of that system (Conant & Ross Ashby, 1970). However, some imagination is needed to take this lesson from the actual theorem, which technically only states that there exists an optimal regulator that is a deterministic function of the state of the system (which could be trivial, Wentworth (2021)). Our theorem less ambiguously states that any robust agent must have learned an (approximate) causal model of the environment, as described in Section 3.2. It can therefore be interpreted as a more precise, causal good regulator theorem.

Emergent capabilities. Causal models enable a kind of general competency—an agent can use a causal model of its environment to identify regret-bounded policies for any objective $U(\mathbf{Pa}_U \subseteq \mathbf{V})$ without additional data (Theorem 3). This suggests that causal world models could help explain how general competence can arise from narrow training objectives (Brown et al., 2020; Silver et al., 2021). By Theorems 1 and 2, agents trained to maximise reward across many environments are incentivized to learn a causal world model (as they cannot generalise without one), which can in turn be used to solve any other decision task in the environment. This incentive does not imply that training an agent with a simple reward signal is sufficient to learn causal world models. E.g. it will still be impossible for an agent to learn a causal model (and therefore to generalise) if the model is not causally identifiable from its training data. The question is then if current methods and training schemes are sufficient for learning causal world models. Early results suggest that transformer models can learn world models capable of out-of-distribution prediction (Li et al., 2022) (see Section 6 for discussion). While foundation models are capable of achieving state of the art accuracy on causal reasoning benchmarks (Kiciman et al., 2023), how they achieve this (and if it constitutes bona fide causal reasoning) is debated (Zečević et al., 2023).

Causal discovery. Theorems 1 and 2 involve learning the causal structure of the environment by observing the agent’s policy under interventions. It is perhaps surprising that the response of this single variable to interventions is sufficient to identify all associative and causal relations in \mathbf{Anc}_U . Typically, causal discovery algorithms involve measuring the response of many variables to interventions (Vowels et al., 2022). Also, many causal discovery algorithms assume independent causal mechanisms (Schölkopf et al., 2021), which is equivalent to assuming no agents are present in the data generating process (Kenton et al., 2023). However, our results suggest that agents could be powerful resources for causal discovery. In Appendix B we use the proof of Theorem 2 to derive a causal discovery algorithm for learning causal structure over latents, and test it on synthetic data.

Applicability of causal methods. Causal models have been used to formally define concepts such as intent (Halpern & Kleiman-Weiner, 2018; Ward et al., 2024), harm (Richens et al., 2022), deception (Ward et al., 2023b), manipulation (Ward et al., 2023a) and incentives (Everitt et al., 2021), and are required for approaches to explainability (Wachter et al., 2017) and fairness (Kusner et al., 2017). Methods for designing safe and ethical AI systems that build on these definitions require causal models of the data generating process, which are typically hard to learn, leading some to doubt their practicality (Fawkes et al., 2022; Rahmattalabi & Xiang, 2022). However, our results show

that learning a causal model of the data generating process is necessary, and demonstrate that these models can be elicited from sufficiently robust agents. These could in turn be used to support causal methods for safety and fairness.

Limitations. Theorems 1 and 2 require agents to be robust to a large set of domain shifts (local interventions on all environment variables). Theorem 2 shows that loosening regret bounds results in some causal relations being unidentifiable from the agent’s policy. Hence, we expect it is still possible to learn some causal knowledge of the environment from agents that are robust to a smaller set of domain shifts, albeit less complete than the full underlying CBN. Finally, our results only apply to unmediated decision tasks (Assumption 1). We expect Theorems 1 and 2 can be extended to active decision tasks, as Assumption 1 does not play a major role beyond simplifying the proofs.

5 RELATED WORK

Several recent empirical works have explored if deep learning models learn ‘surface statistics’ (e.g. correlations between inputs and outputs) or learn internal representations of the world (McGrath et al., 2022; Abdou et al., 2021; Li et al., 2022; Gurnee & Tegmark, 2023). Our results offer some theoretical clarity to this discussion, tying an agent’s performance to the fidelity of its world model, and that going beyond ‘surface statistics’ to learning causal relations is fundamentally necessary for robustness. One study in particular (Li et al., 2022) found that a GPT model trained to predict legal next moves in the board game Othello learned a linear representation of the board state (Nanda, 2023). Further, this internal representation of the board state could be changed by intervening on the intermediate activations, with the model updating its predictions consistent with the intervention, including interventions that take the board state outside of the training distribution. This indicates that the network is learning and utilising a representation of the data generating process that can support out-of-distribution generalisation under interventions—much like a causal model.

The problem of evaluating policies under distributional shifts has been studied extensively in causal transportability (CT) theory (Bareinboim & Pearl, 2016; Bellot & Bareinboim, 2022). CT aims to provide necessary and sufficient conditions for policy evaluation under known domain shifts when all assumptions on the data generating process (i.e. inductive biases) can be expressed as constraints on causal structure (Bareinboim & Pearl, 2016). However, deep learning algorithms can exploit a much larger set of inductive biases (Neyshabur et al., 2014; Battaglia et al., 2018; Rahaman et al., 2019; Goyal & Bengio, 2022) which in many real-world tasks may be sufficient to identify low regret policies without requiring causal knowledge. Thus, CT does not imply that agents must learn causal models in order to generalise unless we assume agents only use causal assumptions to begin with, which would be proof by assumption. See Appendix G for further discussion.

A similar result to Theorems 1 and 2 is the causal hierarchy theorem (CHT) (Bareinboim et al., 2022; Ibeling & Icard, 2021), which shows that observational data is almost always *insufficient* for identifying all causal relations between environment variables, whereas our results state that the set of optimal policies is almost always sufficient to identify all causal relations. In Section 3.2 we discuss the similarities between these theorems, and in Appendix G we discuss their differences.

6 CONCLUSION

Causal reasoning is foundational to human intelligence, and has been conjectured to be necessary for achieving human level AI (Pearl, 2019). In recent years, this conjecture has been challenged by the development of artificial agents capable of generalising to new tasks and domains without explicitly learning or reasoning on a causal model. And while the necessity of causal models for solving causal inference tasks has been established (Bareinboim et al., 2022), their role in decision tasks such as classification and reinforcement learning is less clear.

We have resolved this conjecture in a model-independent way, showing that any agent capable of robustly solving a decision task must have learned a causal model of the data generating process, regardless of how the agent is trained or the details of its architecture. This hints at an even deeper link between causality and general intelligence, as this causal model can be used to simulate the environment and do policy evaluation for any given objective function. By establishing this formal connection between causality and generalisation, our results show that causal world models are a necessary ingredient for robust and general AI.

Acknowledgements. We would like to thank Alexis Bellot, Damiano Fornasiere, Pietro Greiner, James Fox, Matt MacDermott, David Reber, David Watson and Philip Bachman for their helpful discussions and comments on the manuscript.

REFERENCES

- Mostafa Abdou, Artur Kulmizev, Daniel Hershcovich, Stella Frank, Ellie Pavlick, and Anders Søgaard. Can language models encode perceptual structure without grounding? a case study in color. *arXiv preprint arXiv:2109.06129*, 2021.
- Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, 2019.
- Elias Bareinboim and Judea Pearl. Controlling selection bias in causal inference. In *Artificial Intelligence and Statistics*, pp. 100–108. PMLR, 2012a.
- Elias Bareinboim and Judea Pearl. Transportability of causal effects: Completeness results. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 26, pp. 698–704, 2012b.
- Elias Bareinboim and Judea Pearl. Causal inference and the data-fusion problem. *Proceedings of the National Academy of Sciences*, 113(27):7345–7352, 2016.
- Elias Bareinboim, Juan D Correa, Duligur Ibeling, and Thomas Icard. On pearl’s hierarchy and the foundations of causal inference. In *Probabilistic and Causal Inference: The Works of Judea Pearl*, pp. 507–556. 2022.
- Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.
- Alexis Bellot and Elias Bareinboim. Partial transportability for domain generalization. 2022.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Daniel C Castro, Ian Walker, and Ben Glocker. Causality matters in medical imaging. *Nature Communications*, 11(1):3673, 2020.
- Taco Cohen and Max Welling. Group equivariant convolutional networks. In *International conference on machine learning*, pp. 2990–2999. PMLR, 2016.
- Roger C Conant and W Ross Ashby. Every good regulator of a system must be a model of that system. *International journal of systems science*, 1(2):89–97, 1970.
- Juan Correa and Elias Bareinboim. A calculus for stochastic interventions: Causal effect identification and surrogate experiments. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 10093–10100, 2020.
- A P Dawid. Influence diagrams for causal modelling and inference. *International Statistical Review / Revue Internationale de Statistique*, 70:161–189, 2002.
- Daniel C Dennett. *The intentional stance*. MIT press, 1989.
- Logan Engstrom, Brandon Tran, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry. Exploring the landscape of spatial robustness. In *International conference on machine learning*, pp. 1802–1811. PMLR, 2019.
- Tom Everitt, Ryan Carey, Eric D Langlois, Pedro A Ortega, and Shane Legg. Agent incentives: A causal perspective. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 11487–11495, 2021.

- Abolfazl Farahani, Sahar Voghoei, Khaled Rasheed, and Hamid R Arabnia. A brief review of domain adaptation. *Advances in Data Science and Information Engineering: Proceedings from ICDATA 2020 and IKE 2020*, pp. 877–894, 2021.
- Jake Fawkes, Robin Evans, and Dino Sejdinovic. Selection, ignorability and challenges with causal fairness. In *Conference on Causal Learning and Reasoning*, pp. 275–289. PMLR, 2022.
- Alison Gopnik, Laura Schulz, and Laura Elizabeth Schulz. *Causal learning: Psychology, philosophy, and computation*. Oxford University Press, 2007.
- Anirudh Goyal and Yoshua Bengio. Inductive biases for deep learning of higher-level cognition. *Proceedings of the Royal Society A*, 478(2266):20210068, 2022.
- Sharut Gupta, Stefanie Jegelka, David Lopez-Paz, and Kartik Ahuja. Context is environment, 2023.
- Wes Gurnee and Max Tegmark. Language models represent space and time. *arXiv preprint arXiv:2310.02207*, 2023.
- Joseph Halpern and Max Kleiman-Weiner. Towards formal definitions of blameworthiness, intention, and moral responsibility. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019.
- Ronald A Howard and James E Matheson. Influence diagrams. *Decision Analysis*, 2(3):127–143, 2005.
- Patrik Hoyer, Dominik Janzing, Joris M Mooij, Jonas Peters, and Bernhard Schölkopf. Nonlinear causal discovery with additive noise models. *Advances in neural information processing systems*, 21, 2008.
- Duligur Ibeling and Thomas Icard. A topological perspective on causal inference. *Advances in Neural Information Processing Systems*, 34:5608–5619, 2021.
- Zachary Kenton, Ramana Kumar, Sebastian Farquhar, Jonathan Richens, Matt MacDermott, and Tom Everitt. Discovering agents. *Artificial Intelligence*, pp. 103963, 2023.
- Emre Kıcıman, Robert Ness, Amit Sharma, and Chenhao Tan. Causal reasoning and large language models: Opening a new frontier for causality. *arXiv preprint arXiv:2305.00050*, 2023.
- Daphne Koller and Brian Milch. Multi-agent influence diagrams for representing and solving games. *Games and economic behavior*, 45(1):181–221, 2003.
- Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. *Advances in neural information processing systems*, 30, 2017.
- Shane Legg and Marcus Hutter. Universal intelligence: A definition of machine intelligence. *Minds and machines*, 17:391–444, 2007.
- Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Emergent world representations: Exploring a sequence model trained on a synthetic task. *arXiv preprint arXiv:2210.13382*, 2022.
- Thomas McGrath, Andrei Kapishnikov, Nenad Tomašev, Adam Pearce, Martin Wattenberg, Demis Hassabis, Been Kim, Ulrich Paquet, and Vladimir Kramnik. Acquisition of chess knowledge in AlphaZero. *Proceedings of the National Academy of Sciences*, 119(47):e2206625119, 2022.
- Christopher Meek. Strong completeness and faithfulness in bayesian networks. *arXiv preprint arXiv:1302.4973*, 2013.
- Nicolai Meinshausen. Causality from a distributional robustness point of view. In *2018 IEEE Data Science Workshop (DSW)*, pp. 6–10. IEEE, 2018.

- Jovana Mitrovic, Dino Sejdinovic, and Yee Whye Teh. Causal inference via kernel deviance measures. *Advances in neural information processing systems*, 31, 2018.
- Joris M Mooij, Jonas Peters, Dominik Janzing, Jakob Zscheischler, and Bernhard Schölkopf. Distinguishing cause from effect using observational data: methods and benchmarks. *The Journal of Machine Learning Research*, 17(1):1103–1204, 2016.
- Neel Nanda. Actually, othello-gpt has a linear emergent world model, Mar 2023. URL <https://neelnanda.io/mechanistic-interpretability/othello>.
- Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. *arXiv preprint arXiv:1412.6614*, 2014.
- Masashi Okamoto. Distinctness of the eigenvalues of a quadratic form in a multivariate sample. *The Annals of Statistics*, pp. 763–765, 1973.
- Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, 2 edition edition, 2009. ISBN 9780521895606.
- Judea Pearl. Theoretical impediments to machine learning with seven sparks from the causal revolution. *arXiv preprint arXiv:1801.04016*, 2018.
- Judea Pearl. The seven tools of causal inference, with reflections on machine learning. *Communications of the ACM*, 62(3):54–60, 2019.
- Judea Pearl and Elias Bareinboim. Transportability of causal and statistical relations: A formal approach. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 25, pp. 247–254, 2011.
- Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. In *2018 IEEE international conference on robotics and automation (ICRA)*, pp. 3803–3810. IEEE, 2018.
- Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *International Conference on Machine Learning*, pp. 5301–5310. PMLR, 2019.
- Aida Rahmattalabi and Alice Xiang. Promises and challenges of causality for ethical machine learning. *arXiv preprint arXiv:2201.10683*, 2022.
- Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al. A generalist agent. *arXiv preprint arXiv:2205.06175*, 2022.
- Hans Reichenbach. *The direction of time*, volume 65. Univ of California Press, 1956.
- Jonathan Richens, Rory Beard, and Daniel H Thompson. Counterfactual harm. *Advances in Neural Information Processing Systems*, 35:36350–36365, 2022.
- Donald B Rubin. Causal inference using potential outcomes: Design, modeling, decisions. *Journal of the American Statistical Association*, 100(469):322–331, 2005.
- Bernhard Schölkopf, Dominik Janzing, Jonas Peters, Eleni Sgouritsa, Kun Zhang, and Joris Mooij. On causal and anticausal learning. *arXiv preprint arXiv:1206.6471*, 2012.
- Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. *Proceedings of the IEEE*, 109(5):612–634, 2021.
- Rohin Shah, Vikrant Varma, Ramana Kumar, Mary Phuong, Victoria Krakovna, Jonathan Uesato, and Zac Kenton. Goal misgeneralization: Why correct specifications aren’t enough for correct goals. *arXiv preprint arXiv:2210.01790*, 2022.

- Zheyang Shen, Peng Cui, Kun Kuang, Bo Li, and Peixuan Chen. Causally regularized learning with agnostic data selection bias. In *Proceedings of the 26th ACM international conference on Multimedia*, pp. 411–419, 2018.
- David Silver, Satinder Singh, Doina Precup, and Richard S Sutton. Reward is enough. *Artificial Intelligence*, 299:103535, 2021.
- Steven A Sloman and David Lagnado. Causality in thought. *Annual review of psychology*, 66: 223–247, 2015.
- Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. *Causation, prediction, and search*. MIT press, 2000.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Matthew J Vowels, Necati Cihan Camgoz, and Richard Bowden. D’ya like dags? a survey on structure learning and causal discovery. *ACM Computing Surveys*, 55(4):1–36, 2022.
- Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without opening the black box: Automated decisions and the GDPR. *Harv. JL & Tech.*, 31:841, 2017.
- Jindong Wang, Cuiling Lan, Chang Liu, Yidong Ouyang, Tao Qin, Wang Lu, Yiqiang Chen, Wenjun Zeng, and Philip Yu. Generalizing to unseen domains: A survey on domain generalization. *IEEE Transactions on Knowledge and Data Engineering*, 2022.
- Francis Rhys Ward, Tom Everitt, Francesco Belardinelli, and Francesca Toni. Honesty is the best policy: defining and mitigating AI deception. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023a.
- Francis Rhys Ward, Francesca Toni, and Francesco Belardinelli. Defining deception in structural causal games. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, pp. 2902–2904, 2023b.
- Francis Rhys Ward, Matt MacDermott, Francesco Belardinelli, Francesca Toni, and Tom Everitt. The reasons that agents act: Intention and instrumental goals. *arXiv preprint arXiv:2402.07221*, 2024.
- John Wentworth. Fixing the good regulator theorem. <https://www.alignmentforum.org/posts/Dx9LoqsEh3gHNJMDk/fixing-the-good-regulator-theorem>, 2021. Accessed: 2023-10-17.
- Garrett Wilson and Diane J Cook. A survey of unsupervised deep domain adaptation. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 11(5):1–46, 2020.
- Yongqin Xian, Christoph H Lampert, Bernt Schiele, and Zeynep Akata. Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly. *IEEE transactions on pattern analysis and machine intelligence*, 41(9):2251–2265, 2018.
- Matej Zečević, Moritz Willig, Devendra Singh Dhami, and Kristian Kersting. Causal parrots: Large language models may talk causality but are not causal. *arXiv preprint arXiv:2308.13067*, 2023.
- Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1):43–76, 2020.