# Amortized Causal Discovery: Learning to Infer Causal Graphs from Time-Series Data

**Sindy Löwe**[*]
UvA-Bosch Delta Lab
University of Amsterdam
loewe.sindy@gmail.com

**David Madras**[*]
University of Toronto
Vector Institute
madras@cs.toronto.edu

**Richard Zemel**
University of Toronto
Vector Institute
CIFAR
zemel@cs.toronto.edu

**Max Welling**
UvA-Bosch Delta Lab
University of Amsterdam
CIFAR
welling.max@gmail.com

## Abstract

Standard causal discovery methods must fit a new model whenever they encounter samples from a new underlying causal graph. However, these samples often share relevant information – for instance, the dynamics describing the effects of causal relations – which is lost when following this approach. We propose Amortized Causal Discovery, a novel framework that leverages such shared dynamics to learn to infer causal relations from time-series data. This enables us to train a single, amortized model that infers causal relations across samples with different underlying causal graphs, and thus makes use of the information that is shared. We demonstrate experimentally that this approach, implemented as a variational model, leads to significant improvements in causal discovery performance, and show how it can be extended to perform well under hidden confounding.

## 1 Introduction

Inferring causal relations in observational time-series is central to many fields of scientific inquiry [7, 48]. Suppose you want to analyze fMRI data, which measures the activity of different brain regions over time — how can you infer the (causal) influence of one brain region on another? This question of inferring causal relations from observational data is addressed by the field of *causal discovery* [16]. While some methods rely on interventions (e.g. randomized trials), performing these is often infeasible, unethical, or too expensive.

In time-series, the assumption that causes temporally precede their effects enables us to discover causal relations in observational data [40], with approaches relying on conditional independence tests [12], heuristic scores [8], or deep learning [49]. All of these methods assume that samples share a single underlying causal graph and refit a new model whenever this assumption does not hold. However, samples with different underlying causal graphs may share relevant information such as the dynamics describing the effects of causal relations. fMRI test subjects may have varying brain connectivity but the same underlying neurochemistry; social networks may have differing structure but comparable interpersonal relationships; different stocks may relate differently to one another but obey similar market forces. Despite a range of relevant applications, inferring causal relations across samples with different underlying causal graphs is as of yet largely unexplored.
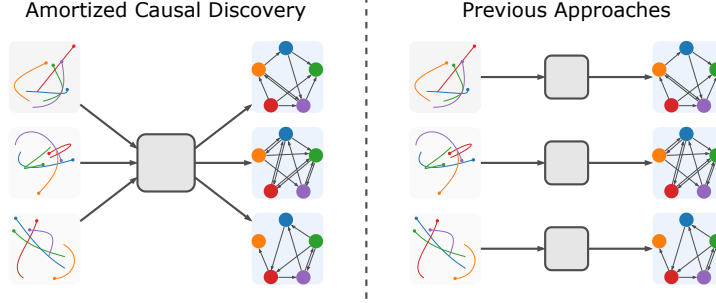
---

[*]equal contribution

Figure 1: Amortized Causal Discovery. We propose to train a single model that infers causal relations across samples with different underlying causal graphs but shared dynamics. This allows us to generalize across samples and to improve our performance with additional training data. In contrast, previous approaches fit a new model for every sample with a different underlying causal graph.

In this paper, we propose a novel causal discovery framework for time-series that embraces this aspect: Amortized Causal Discovery. In this framework, we learn to infer causal relations across samples with different underlying causal graphs but shared dynamics. We achieve this by separating the causal relation prediction from the modeling of their dynamics: an amortized encoder predicts the edges in the causal graph, and a decoder models the dynamics of the system under the predicted causal relations. This setup allows us to pool statistical strength across samples and to achieve significant improvements in performance with additional training data. It also enables "one-shot" inference of causal relations in previously unseen samples without refitting our model. Additionally, we show that Amortized Causal Discovery allows us to improve robustness under hidden confounding, by predicting the values of unobserved variables with the amortized encoder. Our contributions are as follows:

- We formalize Amortized Causal Discovery (ACD), a novel framework for causal discovery in time-series, in which we learn to infer causal relations from samples with different underlying causal graphs but shared dynamics.

- We propose a variational model for ACD, applicable to multi-variate, non-linear data.

- We present experiments demonstrating the effectiveness of this model on a range of causal discovery datasets, both in the fully observed setting and under hidden confounding.

## 2 Background

### 2.1 Granger Causality

Granger causality [18] is one of the most commonly used approaches to infer causal relations from observational time-series data. Its central assumption is that causes precede their effects: if the prediction of the future of time-series $Y$ can be improved by knowing past elements of time-series $X$, then $X$ "Granger causes" $Y$. Originally, Granger causality was defined for linear relations; we follow the more recent definition of Tank et al. [49] for non-linear Granger causality:

**Definition 2.1.** *Non-Linear Granger Causality*: Given N stationary time-series $\boldsymbol{x} = \{\boldsymbol{x}_1, ... \boldsymbol{x}_N\}$ across time-steps $t = \{1, ..., T\}$ and a non-linear autoregressive function $g_j$, such that

$$\boldsymbol{x}_j^{t+1} = g_j(\boldsymbol{x}_1^{\leq t}, ..., \boldsymbol{x}_N^{\leq t}) + \boldsymbol{\varepsilon}_j^t \quad , \tag{1}$$

where $\boldsymbol{x}_j^{\leq t} = (..., \boldsymbol{x}_j^{t-1}, \boldsymbol{x}_j^t)$ denotes the present and past of series $j$ and $\boldsymbol{\varepsilon}_j^t$ represents independent noise. In this setup, time-series $i$ Granger causes $j$, if $g_j$ is not invariant to $\boldsymbol{x}_i^{\leq t}$, i.e. if $\exists \boldsymbol{x}_i'^{\leq t} \neq \boldsymbol{x}_i^{\leq t} : g_j(\boldsymbol{x}_1^{\leq t}, ..., \boldsymbol{x}_i'^{\leq t}, ..., \boldsymbol{x}_N^{\leq t}) \neq g_j(\boldsymbol{x}_1^{\leq t}, ..., \boldsymbol{x}_i^{\leq t}, ... \boldsymbol{x}_N^{\leq t})$.

Granger causal relations are equivalent to causal relations in the underlying directed acyclic graph (DAG) if all relevant variables are observed and no instantaneous[2] connections exist [40].

---

[2]i.e. there are no edges between two variables at the same time step

## 2.2 Neural Relational Inference

Our proposed variational formulation in Section 3.1 is based on the Neural Relation Inference (NRI) model by Kipf et al. [27]. This model infers general, i.e. correlational relations between time-series. For example, based on shared motion patterns between particles (represented as time-series), NRI infers how these particles relate to one another – two particles may be connected by springs while others are connected by static bars – by predicting various edges types in the relational graph. NRI follows the Variational Auto-Encoder framework [25, 42], maximizing the variational lower bound:

$$\mathcal{L} = \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}[\log p_\theta(\boldsymbol{x}|\boldsymbol{z})] - \mathrm{KL}[q_\phi(\boldsymbol{z}|\boldsymbol{x})||p(\boldsymbol{z})] \quad , \tag{2}$$

where the encoder $q_\phi(\boldsymbol{z}|\boldsymbol{x})$ predicts the discrete relation types between time-series, and the decoder $p_\theta(\boldsymbol{x}|\boldsymbol{z})$ models these relation's dynamics to predict the future trajectory of each time-series.

NRI's input samples consist of N time-series $\boldsymbol{x} = \{\boldsymbol{x}_1, ...\boldsymbol{x}_N\}$ across time-steps $t = \{1, ..., T\}$. The encoder is a graph neural network (GNN) [15, 26, 30, 44] that takes these time-series as input and returns a factorized distribution $q_\phi(\boldsymbol{z}|\boldsymbol{x}) = \prod q_\phi(\boldsymbol{z}_{ij}|\boldsymbol{x})$. Its output $\boldsymbol{z}_{ij}$ is modeled as a discrete categorical variable and is a one-hot representation of the edges between time-series: $z_{ij,e} = 1$ expresses that there is a directed edge of type $e$ from time-series $i$ to $j$. The corresponding relational graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ is a complete graph, with vertices $v_i \in \mathcal{V}$ for each time-series $i$, and each $(v_i, v_j)$ connected by an edge of the inferred type.

To backpropagate through the discrete samples of $q_\phi(\boldsymbol{z}|\boldsymbol{x})$, they are relaxed using the Gumbel-Softmax distribution [22, 33]. The resulting sample $\boldsymbol{z}_{ij}$ is input to the decoder together with the feature vectors of the time-series at the current time-step $t$, $\boldsymbol{x}^t = \{\boldsymbol{x}_1^t, ...\boldsymbol{x}_N^t\}$. The decoder is another GNN that models $p_\theta(\boldsymbol{x}|\boldsymbol{z}) = \prod_{t=1}^T p_\theta(\boldsymbol{x}^{t+1}|\boldsymbol{x}^{\leq t}, \boldsymbol{z})$. First, it propagates information across the inferred edges $e$ of the relational graph using separate neural networks $f_e$ for each edge-type:

$$h_{ij}^t = \sum_e z_{ij,e} f_e([\boldsymbol{x}_i^t, \boldsymbol{x}_j^t]) \quad . \tag{3}$$

The decoder then accumulates the incoming messages to each node to predict the change between the current and the next time-step:

$$\boldsymbol{\mu}_j^{t+1} = \boldsymbol{x}_j^t + f_v\left(\left[\sum_{i \neq j} h_{ij}^t, \boldsymbol{x}_j^t\right]\right) \tag{4}$$

$$p_\theta(\boldsymbol{x}_j^{t+1}|\boldsymbol{x}^t, \boldsymbol{z}) = \mathcal{N}(\boldsymbol{\mu}_j^{t+1}, \sigma^2 \mathbb{I}) \quad . \tag{5}$$

In other words, the decoder predicts $\Delta \hat{\boldsymbol{x}}^t$, which is added to the current value of the time-series to yield the prediction for the next time-step $\hat{\boldsymbol{x}}^{t+1} = \boldsymbol{x}^t + \Delta \hat{\boldsymbol{x}}^t$.

# 3 Amortized Causal Discovery

We propose Amortized Causal Discovery (ACD), a framework in which we learn to infer causal relations across samples with different underlying causal graphs but shared dynamics. To illustrate: Suppose you want to infer synaptic connections (i.e. causal relations) between neurons based on their spiking behaviour. You are given a set of $N$ recordings (i.e. samples), each containing $S$ time-series representing the firing of $S$ individual neurons. Even though the wiring of neurons may differ between recordings, the dynamics of how neurons connected by synapses influence one another stays the same. ACD takes advantage of such shared dynamics to improve the prediction of causal relations. In this section, we outline this framework, and describe a probabilistic implementation based on a causal formulation of NRI. We also extend our approach to model unobserved confounders.

**Preliminaries**  We begin with a dataset $\boldsymbol{X} = \{\boldsymbol{x}_s\}_{s=1}^S$ of $S$ samples, where each sample $\boldsymbol{x}_s$ consists of $N$ stationary time-series $\boldsymbol{x}_s = \{\boldsymbol{x}_{s,1}, \ldots, \boldsymbol{x}_{s,N}\}$ across timesteps $t = \{1, ..., T\}$. We denote the $t$-th time-step of the $i$-th time-series of $\boldsymbol{x}_s$ as $\boldsymbol{x}_{s,i}^t$ (sometimes omitting $s$ for brevity). We assume there is a directed acyclic graph $\mathcal{G}_s^{1:T} = \{\mathcal{V}_s^{1:T}, \mathcal{E}_s^{1:T}\}$ underlying the generative process of each sample. This graph is a structural causal model (SCM) [38]. Its endogenous (observed) variables are vertices $v_{s,i}^t \in \mathcal{V}_s^{1:T}$ for each time-series $i$ and each time-step $t$, and every set of incoming edges to an

endogenous variable defines inputs to a deterministic function $g_{s,i}^t$ generating that variable.[3] The edges are defined by ordered pairs of vertices $\mathcal{E}_s^{1:T} = \{(v_{s,i}^t, v_{s,j}^{t'})\}$, which we make two assumptions about:

1. No edges are instantaneous ($t = t'$) or go back in time. Thus, $t < t'$ for all edges.

2. Edges are invariant to time. Thus, if $(v_{s,i}^t, v_{s,j}^{t+k}) \in \mathcal{E}_s^{1:T}$, then $\forall 1 \leq t' \leq T - k : (v_{s,i}^{t'}, v_{s,j}^{t'+k}) \in \mathcal{E}_s^{1:T}$. The associated structural equations $g_{s,i}^t$ are invariant to time as well, i.e. $g_{s,i}^t = g_{s,i}^{t'} \ \forall t, t'$.

The first assumption states that causes temporally precede their effects and makes causal relations identifiable from observational data [39]. The second simplifies modeling: it is a fairly general assumption which allows us to define dynamics that govern all time-steps (Eq. (6)).

Throughout this paper, we are interested in discovering the *summary graph* $\mathcal{G}_s = \{\mathcal{V}_s, \mathcal{E}_s\}$ [40]. It consists of vertices $v_{s,i} \in \mathcal{V}_s$ for each time-series $i$ in sample $s$, and has edges whenever they exist in $\mathcal{E}_s^{1:T}$ at any time-step, i.e. $\mathcal{E}_s = \{(v_{s,i}, v_{s,j}) \mid \exists t, t' : (v_{s,i}^t, v_{s,j}^{t'}) \in \mathcal{E}_s^{1:T}\}$. Note that $\mathcal{G}_s$ is directed but may be cyclic.

**Amortized Causal Discovery**   The key assumption for Amortized Causal Discovery is that there exists some fixed function $g$ that describes the dynamics of *all* samples $\boldsymbol{x}_s \in \boldsymbol{X}$ given their past observations $\boldsymbol{x}_s^{\leq t}$ and their underlying causal graph $\mathcal{G}_s$:

$$\boldsymbol{x}_s^{t+1} = g(\boldsymbol{x}_s^{\leq t}, \mathcal{G}_s) + \boldsymbol{\varepsilon}_s^t \quad . \tag{6}$$

There are two variables in this data-generating process that we would like to model: the dynamics $g$ that are shared across all samples, and the causal graph $\mathcal{G}_s$ that is specific to sample $\boldsymbol{x}_s$. This separation between the dynamics and the causal graph allows us to introduce an amortized causal discovery algorithm $\mathcal{A}$ which learns to infer the causal graph $\mathcal{G}_s$ given the sample $\boldsymbol{x}_s$. Thus, we aim to learn a dynamics model $f$ to approximate $g$, and the causal discovery algorithm $\mathcal{A}$, which are amortized (i.e. shared) across all samples and predict $k \in \mathbb{N}$ steps into the future:

$$\boldsymbol{x}_s^{t+k} \approx f(\boldsymbol{x}_s^{\leq t}, \mathcal{A}(\boldsymbol{x}_s)) \quad . \tag{7}$$

We formalize Amortized Causal Discovery (ACD) as follows. Let $\mathcal{G}_s \in \mathbb{G}$ and any single step, partial or full sequence $\boldsymbol{x}_s^t, \boldsymbol{x}_s^{\leq t}, \boldsymbol{x}_s \in \mathbb{X}$. The model consists of two components: a causal discovery encoder $\mathcal{A} : \mathbb{X} \rightarrow \mathbb{G}$ which infers a graph for each input sample, and a decoder $f : \mathbb{X} \times \mathbb{G} \rightarrow \mathbb{X}$ which models the dynamics. This model is optimized with a sample-wise loss $\ell : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$ which scores how well the decoder models the true dynamics of $\boldsymbol{x}_s$, and a regularization term $r : \mathbb{G} \rightarrow \mathbb{R}$ on the inferred graphs (e.g. enforcing sparsity). Then, given some dataset $\boldsymbol{X}_{train}$ with $S$ samples, we optimize:

$$f_\star, \mathcal{A}_\star = \operatorname{argmin}_{f, \mathcal{A}} \mathcal{L}(f, \mathcal{A}, \boldsymbol{X}_{train}) \tag{8}$$

$$\text{where } \mathcal{L}(f, \mathcal{A}, \boldsymbol{X}_{train}) = \sum_{s=1}^{S} \sum_{t=1}^{T-k} \ell(\boldsymbol{x}_s^{t+k}, f(\boldsymbol{x}_s^{\leq t}, \mathcal{A}(\boldsymbol{x}_s))) + r(\mathcal{A}(\boldsymbol{x}_s)) \quad . \tag{9}$$

Once we have completed optimization, we can perform causal graph prediction on any new input test sample $\boldsymbol{x}_{test}$ in two ways – by using the amortized encoder:

$$\hat{\mathcal{G}}_{Enc} = \mathcal{A}_\star(\boldsymbol{x}_{test}) \quad ; \tag{10}$$

or by optimizing through the learned decoder, which we term *Test-Time Adaptation (TTA)*:

$$\hat{\mathcal{G}}_{TTA} = \operatorname{argmin}_{\mathcal{G}} \mathcal{L}(f_\star, \mathcal{G}, \boldsymbol{x}_{test}) \quad . \tag{11}$$

We can also combine these, e.g., by performing TTA on an initial graph estimate produced by the amortized encoder. As with $\mathcal{G}_s$, $\hat{\mathcal{G}}$ may be cyclic.

By separating the prediction of causal relations from the modeling of their dynamics, ACD yields a number of benefits. ACD can learn to infer causal relations across samples with different underlying causal graphs. By generalizing across samples, it can improve causal discovery performance with increasing training data size. ACD can infer causal relations in previously unseen test-cases in "one-shot", without re-fitting. We can replace either $f$ or $\mathcal{A}$ with ground truth annotations, or simulate the outcome of counterfactual causal relations. Additionally, ACD can be applied in the standard causal discovery setting, where only a single causal graph underlies all samples, by replacing the amortized encoder $\mathcal{A}$ with an estimated graph $\hat{\mathcal{G}}$ (or distribution over $\mathbb{G}$) in Eq. (8).

---

[3]The SCM also includes an exogenous (unobserved), independently-sampled error variable $\epsilon_v$ as a parent of each vertex $v$, which we do not model and thus leave out for brevity.
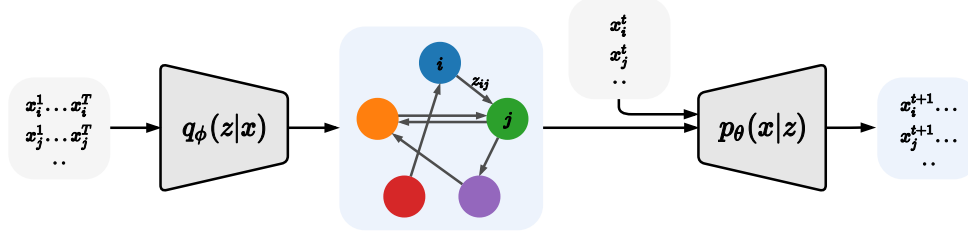
Figure 2: ACD-NRI: A Probabilistic Approach to Amortized Causal Discovery. An amortized encoder $q_\phi(z|x)$ predicts the causal relations between the input time-series $x$. A decoder $p_\theta(x|z)$ models the dynamics of the future trajectories $x^{t+1}$ given their current values $x^t$ and the predicted relations $z$. This separation between causal relation prediction and modeling lets us train the model across samples with different underlying causal graphs but shared dynamics.

## 3.1 A Probabilistic Implementation of ACD

We take a probabilistic approach to ACD and model the functions $f$ and $\mathcal{A}$ using variational inference. We amortize the encoder with a function $q_\phi(z|x)$, which outputs a distribution over $z$ representing the predicted edges $\hat{\mathcal{E}}$ in the causal graph; and we learn a decoder $p_\theta(x|z)$ which probabilistically models the dynamics of the time-series. We choose a negative log-likelihood for the decoder loss $\ell$ and a KL-Divergence to a prior distribution over $\mathbb{G}$ for the regularizer $r$. As a result, our loss function $\mathcal{L}$ is a variational lower bound as in Eq. (2).

Following the variational formulation of ACD, we can adopt a wide range of encoder/decoder schemes to instantiate the model. We choose the NRI model [27] (described in Section 2.2), which uses a graph neural network encoder and decoder. It predicts an edge type $e$ for each pair of nodes $(v_i, v_j)$, and models it with some learned function $f_e$. In order to align it with the philosophy of Granger Causality, we include a "no edge"-type edge function along which no information is propagated.[4] This yields the model we call ACD-NRI, where Eq. (3) becomes

$$h_{ij}^t = \begin{cases} 0 & \text{if } z_{ij,0} = 1 \\ \sum_e z_{ij,e} f_e([x_i^t, x_j^t]) & \text{else} \end{cases} . \tag{12}$$

Thus, if the encoder predicts the "no edge"-type edge $e = 0$ for $z_{ij}$ by setting $z_{ij,0} = 1$, the decoder uses the zero function and no information is propagated from time-series $i$ to $j$. Due to this, we can draw a connection to Granger causality: if $z_{ij,0} = 1$, time series $i$ does not Granger cause the decoder's prediction for $j$ (see Appendix A). As a result, we expect ACD-NRI to predict "no edge" where no (Granger) causal connection exists, and thus to estimate the true causal graph [40].

## 3.2 Hidden Confounding

When hidden confounders exist, Granger causality is not guaranteed to correspond to the true causal graph [40]. For instance, if an unobserved time-series $U$ causes both time-series $X$ and $Y$, then the past of $X$ can help predict the future of $Y$, even though there is no causal link between them. This can raise large practical roadblocks to causal discovery.

Some empirical work shows that encoder-based models with enough proxies (variables caused by hidden confounders) can improve causal inference under hidden confounding [32, 36], and theoretical work proves the identifiability of latent variables from proxies under some assumptions [2, 28]. Inspired by this, we extend the amortized encoder $q_\phi(z|x)$ to additionally predict the values of the unobserved variables. We encourage it to represent the unobserved variable by implementing a structural bias – depending on the type of unobserved variable that we model, its predicted value is fed into the remaining model differently. The decoder remains responsible for modeling the dynamics, and now also processes the predictions for the unobserved variable. For details, see Section 5.2.

---

[4]As shown by Peters et al. [40], time-series $i$ does not Granger cause $j$ in the observed data if there is no edge $(v_i, v_j)$ in the underlying graph, assuming no instantaneous connections and enough samples.
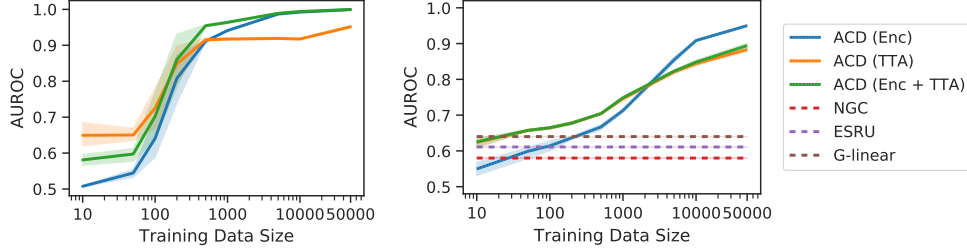
Figure 3: Influence of training data size on causal discovery performance (in AUROC) on the particles dataset (A-left) and Kuramoto (B-right). ACD-NRI improves with additional training data (*Enc*), with test-time adaptation (*TTA* and *Enc+TTA*) boosting results in low-data settings.

## 4 Related Work

A range of approaches to causal discovery in both temporal and non-temporal data exist [19, 40, 48]. One common class is *constraint-based*, relying on conditional independence testing to uncover an underlying DAG structure or equivalence class [48]. These methods predict a single graph $\hat{\mathcal{G}}$ (or equivalence class) for all samples. There is no notion of fitting a dynamics model for time-series methods in this class [12]. Another common class of methods for causal discovery is *score-based* [8]. Here, a heuristic score function $h$ is chosen, and the methods perform a search through graph space to optimize this score, i.e. $\hat{\mathcal{G}} = \operatorname{argmin}_{\mathcal{G}} h(\mathcal{G})$. Our proposed decoder-based inference (Eq. (11)) can be seen as score-based causal discovery with a *learned* score function $\mathcal{L} \circ f_{\star}$.

A third class of methods fits a (possibly regularized) dynamics model $f$ and then analyzes its form to produce a causal graph estimate. This approach can be formulated as $\hat{\mathcal{G}} = \mathcal{A}(\operatorname{argmin}_{f} \mathcal{L}(f, \boldsymbol{X}))$, where $\mathcal{A}$ is the algorithm analyzing the optimized function $f_{\star}$ and $\mathcal{L}$ is some appropriate loss. Note, that these methods do *not* separate causal discovery from dynamics modeling as we propose, and as a result, need to fit a new model whenever they encounter samples with a different underlying causal graph. This class includes vector autoregressive models [21] and a range of recent works that make use of deep learning. Tank et al. [49] and Khanna and Tan [23] use deep recurrent architectures to model $f$ and enforce sparsity on a specific set of weights during training. Nauta et al. [35] follow a similar idea using attention modules. Wu et al. [52] infer causal relations by evaluating how their trajectory prediction is altered by noise injected on the input. Other approaches to causal discovery in temporal data use independence or additivity assumptions [10, 39], Granger causality [3, 4, 11], hidden confounding [14, 34], or seek to predict the temporal direction [6, 41, 46].

Several works have used graph neural networks [5, 27, 43] or attention mechanisms [13, 17, 50, 51] to infer relations between time-series. Alet et al. [1] propose a meta-learning algorithm to additionally model unobserved variables. While these approaches model object relations in a number of ways, they are not explicitly designed to infer *causal* graphical structure.

## 5 Experiments

**Implementation**   In our experiments, the encoder consists of fully-connected networks (MLPs) or 1D CNNs with attentive pooling [31] and the decoder consists of an MLP. They implement two and one edge-propagation steps along the graph, respectively. We measure causal discovery performance by area under the receiver operator curve (AUROC) of predicted edge probabilities over test samples. We compare to recurrent models (Khanna and Tan [23], Tank et al. [49]), and a mutual-information (MI) based model by Wu et al. [52] and several baselines implemented by those authors, including MI (unmodified), transfer entropy [45], and a linear Granger causality. More details in Appendix B; our code is available at `github.com/loeweX/AmortizedCausalDiscovery`.

### 5.1   Fully Observed Amortized Causal Discovery

We test ACD-NRI on two fully-observed physics simulations following non-linear causal dynamics: a simulation of particles exerting spring-like forces on each other; and a simulation of phase-coupled

oscillators (Kuramoto model) [29]. For both, we generate 50,000 training and 10,000 validation samples. We restrict the number of test samples to 200, since the baselines we compare to must be refit for each individual sample. Note, in contrast to the datasets used in Kipf et al. [27], we allow connectivity matrices to be asymmetric to represent causal links in the summary graph of each sample.

**Particles** In our first experiment, we model five particles that move around a two-dimensional space, with some particles influencing others uni-directionally by pulling them with a spring. We find that ACD-NRI learns to infer the causal relations of the test samples almost perfectly, with **0.999 AUROC**. We do not evaluate our baselines on this dataset since they are intended for one-dimensional time-series. Fig. 3A shows that ACD-NRI takes advantage of additional training data to achieve this strong performance (*Enc*). We can further improve our method's performance by replacing the amortized encoder $q_\phi(z|x)$ with a test-time adapted (TTA), non-amortized, variational distribution $q(z)$ that is optimized through the frozen decoder (Eq. (11)). On this dataset, we find that performing TTA on a random $q(z)$ (*TTA*) performs best in the low-data setting, while initializing with the trained encoder (*Enc+TTA*) improves with more data.

**Kuramoto** Next, we test our method on the Kuramoto dataset, which contains 1-D time-series of phase-coupled oscillators. With 50,000 training samples, our method achieves **0.952 AUROC**, vastly outperforming all of the six baselines we compare against (Table 1). Again, we find that ACD-NRI improves with additional training data (Fig. 3B). Additionally, using test-time adaptation (*TTA* and *Enc+TTA*) allows us to outperform the baseline methods with as few as 50 training samples. This benefit of TTA can be largely attributed to two effects. First, TTA closes the amortization gap of the encoder

| Method | AUROC |
|---|---|
| MPIR [52] | $0.502 \pm 0.006$ |
| Transfer Entropy [45] | $0.560 \pm 0.005$ |
| NGC [49] | $0.574 \pm 0.018$ |
| eSRU [23] | $0.607 \pm 0.001$ |
| Mutual Information | $0.616 \pm 0.000$ |
| Linear Granger Causality | $0.647 \pm 0.003$ |
| Amortized Causal Discovery | $\mathbf{0.952 \pm 0.003}$ |

Table 1: AUROC for causal discovery on Kuramoto dataset. 95% confidence interval shown.

$q_\phi(z|x)$ [9]. Second, $q(z)$ overcomes the encoder's overfitting on the training data (as seen in the training curves in Appendix B.2) by being adapted to the individual test samples. Interestingly, on the Kuramoto dataset, the encoder-seeded *TTA* does not consistently outperform the encoder itself.

**Netsim** We apply ACD-NRI to the Netsim dataset [47] of simulated fMRI data and infer the underlying connectivity between brain regions, which is consistent across the 50 samples. Since one graph underlies all samples, we do not expect to outperform other methods, as we cannot benefit from an amortized encoder. Rather, we use a global latent distribution $q(z)$ and optimize its parameters with a decoder, and then use TTA. Nonetheless, we show that our method performs comparably in this setting to methods that are intended for use in the single-graph setting.

| Method | AUROC |
|---|---|
| MPIR [52] | $0.484 \pm 0.017$ |
| Transfer Entropy [45] | $0.543 \pm 0.003$ |
| NGC [49] | $0.624 \pm 0.020$ |
| eSRU [23] | $0.670 \pm 0.015$ |
| Mutual Information | $\mathbf{0.728 \pm 0.002}$ |
| Linear Granger Causality | $0.503 \pm 0.004$ |
| Amortized Causal Discovery | $\mathbf{0.688 \pm 0.051}$ |

Table 2: AUROC for causal discovery on Netsim dataset. 95% confidence interval shown.

## 5.2 Amortized Causal Discovery with Unobserved Variables

### 5.2.1 Latent Temperature

In this experiment, we use the particles dataset and vary an unobserved temperature variable, which modulates how strongly the particles exert force on each other – higher temperatures result in stronger forces and a more chaotic system. For each $x_s$, we sample an independent temperature $c \sim \text{Categorical}([\frac{\alpha}{2}, \alpha, 2\alpha])$ from a categorical distribution with $\alpha \in \mathbb{R}$ and equal probabilities. We predict this unobserved temperature by extending the amortized encoder with an additional latent variable which models a uniform distribution. Then, we add a KL-Divergence between this posterior and a uniform prior on the interval $[0, 4\alpha]$ to our variational loss. To allow for learning in this setting, we introduce an inductive bias: we use a decoder which matches the true dynamics $g$ given the predicted temperature and causal relations. See Appendix C.1 for more details and additional results.
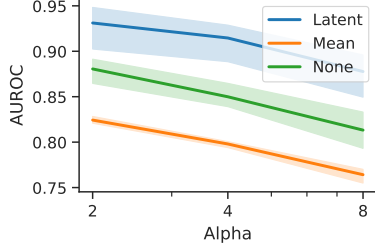
Figure 4: AUROC with unobserved temperature. ACD-NRI with a *latent* variable outperforms a baseline which imputes a *mean* temperature, and a learned fixed-temperature decoder (*None*).
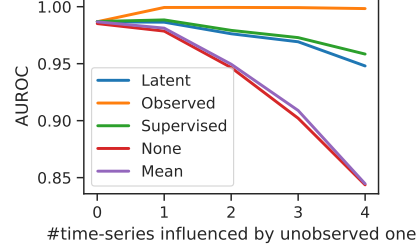
Figure 5: AUROC with unobserved time-series. As more time-series are influenced by the unobserved one (x-axis), the benefit of using an additional *latent* variable for modeling its effects grows.

**Results** Fig. 4 shows the causal discovery results across different values of $\alpha$. ACD-NRI enhanced with an additional latent variable (*Latent*) outperforms both tested baselines across all temperatures: *Mean*, which uses the same ground-truth decoder as *Latent* and fixes the decoder temperature to be the mean of the categorical distribution, and *None*, which does not model $c$ explicitly and trains an MLP decoder. Additionally, this method achieves high predictive performance on the unobserved temperature variable: for $\alpha = 2$, temperature prediction obtains $0.888$ $R^2$, $0.966$ AUROC and $0.644$ accuracy. These results indicate that we can predict and model an unobserved temperature variable, and thus improve robustness under hidden confounding.

### 5.2.2 Unobserved Time-Series

Here, we treat one of the original time-series in the particles dataset as unobserved. This unobserved time-series exhibits the same dynamics as the observed time-series, evolving and causally influencing others the same way as before. We model this unobserved time-series by extending the amortized encoder with an additional latent variable and applying a suitable structural bias: the latent prediction $z_u^t$ for time-steps $t = \{1, ..., T\}$ is treated in the same way as the observed time-series $x$. Its entire trajectory is used by the encoder to predict causal relations, and its value at the current time-step is fed into the decoder. See Appendix C.2 for more details and additional results.
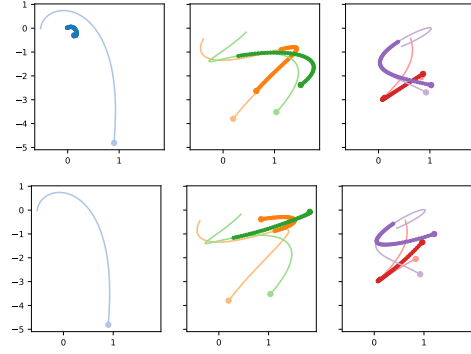


Figure 6: Trajectory prediction with an unobserved time-series (TS). Faded: ground truth. Bold: prediction, starts after observing the first half of the ground truth. Dots denote end of TS. Top: ACD-NRI with *Latent*, bottom: *None* baseline - does not model unobserved TS. Left: unobserved TS, middle: TS directly influenced by unobserved, right: remaining TS. Though we underestimate the unobserved TS, observed TS prediction improves.

**Results** Fig. 5 shows how the causal discovery AUROC depends on the number of observed time-series directly influenced by the unobserved one. When this number is zero, all tested approaches perform the same. With growing numbers of influenced time-series, the baselines that either ignore the missing time-series (*None*) or impute its value with the average of the observed time-series over time (*Mean*) deteriorate strongly. In contrast, the proposed ACD-NRI with a *Latent* variable stays closer to the performance of the fully *Observed* baseline. As shown in Fig. 6, it also improves the future trajectory prediction of the observed time-series. It is only slightly outperformed by a *Supervised* method that optimizes the prediction of the unobserved time-series based on mean-squared error (MSE) with ground-truth trajectories. These results indicate that ACD-NRI can use latent variables to improve robustness to unobserved time-series.

## 6 Conclusion

In this paper, we introduce ACD, a framework for causal discovery in time series data which can leverage the information that is shared across samples. We provide a probabilistic implementation, ACD-NRI, and demonstrate significant performance gains when predicting causal relations. Addi-

tionally, we maintain our performance under hidden confounding. Exciting future directions include interventions, more flexible graph structures, or methods that adapt dynamically to the type of hidden confounder at hand.

## Broader Impact

This paper presents, first and foremost, fundamental research, and as such does not necessarily have direct applications. Additionally, causal discovery as a field is not yet widely applied in any particular area. However, a wide range of future applications inspired by the ideas in this paper are possible, some of which are outlined in the paper. They are mostly in the area of assisting inquiry into observed phenomena. For example, causal discovery algorithms for time series data could help with understanding networks of neural signals and brain activity, for social network analysis, or for modeling of physical systems; in particular, they can improve understanding of such systems while alleviating the need for interventions. These algorithms could be implemented by a variety of actors in pursuit of various incentives, assuming sufficiently large and clean data, or enough domain knowledge.

## Acknowledgments and Disclosure of Funding

## References

[1] F. Alet, E. Weng, T. Lozano-Pérez, and L. P. Kaelbling. Neural relational inference with fast modular meta-learning. In *Advances in Neural Information Processing Systems*, pages 11804–11815, 2019.

[2] E. S. Allman, C. Matias, J. A. Rhodes, et al. Identifiability of parameters in latent structure models with many observed variables. *The Annals of Statistics*, 37(6A):3099–3132, 2009.

[3] M. T. Bahadori and Y. Liu. An examination of practical granger causality inference. In *Proceedings of the 2013 SIAM International Conference on Data Mining*, pages 467–475. SIAM, 2013.

[4] L. Barnett, A. B. Barrett, and A. K. Seth. Granger causality and transfer entropy are equivalent for gaussian variables. *Physical Review Letters*, 103(23):238701, 2009.

[5] P. Battaglia, R. Pascanu, M. Lai, D. J. Rezende, et al. Interaction networks for learning about objects, relations and physics. In *Advances in Neural Information Processing Systems*, pages 4502–4510, 2016.

[6] S. Bauer, B. Schölkopf, and J. Peters. The arrow of time in multivariate time series. In *International Conference on Machine Learning (ICML)*, pages 2043–2051, 2016.

[7] C. Berzuini, P. Dawid, and L. Bernardinell. *Causality: Statistical perspectives and applications*. John Wiley & Sons, 2012.

[8] D. M. Chickering. Optimal structure identification with greedy search. *Journal of Machine Learning Research*, 3(Nov):507–554, 2002.

[9] C. Cremer, X. Li, and D. Duvenaud. Inference suboptimality in variational autoencoders. In *International Conference on Machine Learning (ICML)*, pages 1078–1086, 2018.

[10] M. Eichler. *Causal inference in time series analysis*. Wiley Online Library, 2012.

[11] M. Eichler. Graphical modelling of multivariate time series. *Probability Theory and Related Fields*, 153(1-2):233–268, 2012.

[12] D. Entner and P. O. Hoyer. On causal discovery from time series data using FCI. *Probabilistic Graphical Models*, pages 121–128, 2010.

[13] F. B. Fuchs, A. R. Kosiorek, L. Sun, O. P. Jones, and I. Posner. End-to-end recurrent multi-object tracking and trajectory prediction with relational reasoning. *arXiv preprint arXiv:1907.12887*, 2019.

[14] P. Geiger, K. Zhang, B. Schoelkopf, M. Gong, and D. Janzing. Causal inference by identification of vector autoregressive processes with hidden components. In *International Conference on Machine Learning (ICML)*, pages 1917–1925, 2015.

[15] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. Neural message passing for quantum chemistry. In *International Conference on Machine Learning (ICML)*, pages 1263–1272, 2017.

[16] C. Glymour, K. Zhang, and P. Spirtes. Review of causal discovery methods based on graphical models. *Frontiers in Genetics*, 10, 2019.

[17] A. Goyal, A. Lamb, J. Hoffmann, S. Sodhani, S. Levine, Y. Bengio, and B. Schölkopf. Recurrent independent mechanisms. *arXiv preprint arXiv:1909.10893*, 2019.

[18] C. W. Granger. Investigating causal relations by econometric models and cross-spectral methods. *Econometrica: Journal of the Econometric Society*, pages 424–438, 1969.

[19] C. Heinze-Deml, M. H. Maathuis, and N. Meinshausen. Causal structure learning. *Annual Review of Statistics and Its Application*, 5:371–391, 2018.

[20] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.

[21] A. Hyvärinen, K. Zhang, S. Shimizu, and P. O. Hoyer. Estimation of a structural vector autoregression model using non-gaussianity. *Journal of Machine Learning Research*, 11(56):1709–1731, 2010.

[22] E. Jang, S. Gu, and B. Poole. Categorical reparametrization with gumble-softmax. In *International Conference on Learning Representations (ICLR)*, 2017.

[23] S. Khanna and V. Y. Tan. Economy statistical recurrent units for inferring nonlinear granger causality. *International Conference on Learning Representations (ICLR)*, 2019.

[24] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015.

[25] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

[26] T. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*, 2017.

[27] T. Kipf, E. Fetaya, K.-C. Wang, M. Welling, and R. Zemel. Neural relational inference for interacting systems. In *International Conference on Machine Learning (ICML)*, pages 2688–2697, 2018.

[28] J. B. Kruskal. Three-way arrays: rank and uniqueness of trilinear decompositions, with application to arithmetic complexity and statistics. *Linear Algebra and its Applications*, 18(2):95–138, 1977.

[29] Y. Kuramoto. Self-entrainment of a population of coupled non-linear oscillators. In *International Symposium on Mathematical Problems in Theoretical Physics*, pages 420–422. Springer, 1975.

[30] Y. Li, D. Tarlow, M. Brockschmidt, and R. Zemel. Gated graph sequence neural networks. In *International Conference on Learning Representations (ICLR)*, 2016.

[31] Z. Lin, M. Feng, C. N. d. Santos, M. Yu, B. Xiang, B. Zhou, and Y. Bengio. A structured self-attentive sentence embedding. *arXiv preprint arXiv:1703.03130*, 2017.

[32] C. Louizos, U. Shalit, J. M. Mooij, D. Sontag, R. Zemel, and M. Welling. Causal effect inference with deep latent-variable models. In *Advances in Neural Information Processing Systems*, pages 6446–6456, 2017.

[33] C. J. Maddison, A. Mnih, and Y. W. Teh. The concrete distribution: A continuous relaxation of discrete random variables. In *International Conference on Learning Representations (ICLR)*, 2017.

[34] D. Malinsky and P. Spirtes. Causal structure learning from multivariate time series in settings with unmeasured confounding. In *Proceedings of 2018 ACM SIGKDD Workshop on Causal Discovery*, pages 23–47, 2018.

[35] M. Nauta, D. Bucur, and C. Seifert. Causal discovery with attention-based convolutional neural networks. *Machine Learning and Knowledge Extraction*, 1(1):312–340, 2019.

[36] S. Parbhoo, M. Wieser, A. Wieczorek, and V. Roth. Information bottleneck for estimating treatment effects with systematically missing covariates. *Entropy*, 22(4):389, 2020.

[37] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pages 8024–8035, 2019.

[38] J. Pearl. *Causality*. Cambridge University Press, 2009.

[39] J. Peters, D. Janzing, and B. Schölkopf. Causal inference on time series using restricted structural equation models. In *Advances in Neural Information Processing Systems*, pages 154–162, 2013.

[40] J. Peters, D. Janzing, and B. Schölkopf. *Elements of causal inference: foundations and learning algorithms*. MIT Press, 2017.

[41] L. C. Pickup, Z. Pan, D. Wei, Y. Shih, C. Zhang, A. Zisserman, B. Scholkopf, and W. T. Freeman. Seeing the arrow of time. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2035–2042, 2014.

[42] D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning (ICML)*, pages 1278–1286, 2014.

[43] A. Santoro, D. Raposo, D. G. Barrett, M. Malinowski, R. Pascanu, P. Battaglia, and T. Lillicrap. A simple neural network module for relational reasoning. In *Advances in Neural Information Processing Systems*, pages 4967–4976, 2017.

[44] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2008.

[45] T. Schreiber. Measuring information transfer. *Physical Review Letters*, 85(2):461, 2000.

[46] S. Shimizu, T. Inazumi, Y. Sogawa, A. Hyvärinen, Y. Kawahara, T. Washio, P. O. Hoyer, and K. Bollen. DirectLiNGAM: A direct method for learning a linear non-gaussian structural equation model. *Journal of Machine Learning Research*, 12(Apr):1225–1248, 2011.

[47] S. M. Smith, K. L. Miller, G. Salimi-Khorshidi, M. Webster, C. F. Beckmann, T. E. Nichols, J. D. Ramsey, and M. W. Woolrich. Network modelling methods for fMRI. *Neuroimage*, 54(2):875–891, 2011.

[48] P. Spirtes, C. N. Glymour, R. Scheines, and D. Heckerman. *Causation, prediction, and search*. MIT Press, 2000.

[49] A. Tank, I. Covert, N. Foti, A. Shojaie, and E. Fox. Neural granger causality for nonlinear time series. *arXiv preprint arXiv:1802.05842*, 2018.

[50] S. Van Steenkiste, M. Chang, K. Greff, and J. Schmidhuber. Relational neural expectation maximization: Unsupervised discovery of objects and their interactions. In *International Conference on Learning Representations (ICLR)*, 2018.

[51] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017.

[52] T. Wu, T. Breuel, M. Skuhersky, and J. Kautz. Discovering nonlinear relations with minimum predictive information regularization. *arXiv preprint arXiv:2001.01885*, 2020.

# A  Granger Causality in NRI

Here, we show that when constraining the edge-type $e = 0$ to be the zero function, time series $i$ does not Granger cause the model prediction of $j$ in ACD-NRI.

**Claim**: If $z_{ij,0} = 1$, $i$ does not Granger cause the model prediction of $j$ in ACD-NRI.

*Proof.* According to Definition 2.1, time-series $i$ does not cause $j$, if $g_j$ is invariant to $\boldsymbol{x}_i^{\leq t}$. In our model, the decoder represents this non-linear model $g_j$ and consists of two functions. First, it propagates information across edges using Eq. (12). This function returns a value of zero, if $z_{ij,0} = 1$. This output is used for the second function, described in Eq. (1), which does not introduce any new terms that depend on $i$. Thus, if $z_{ij,0} = 1$, the decoder's prediction for $j$ is invariant to $\boldsymbol{x}_i^{\leq t}$, and $i$ does not Granger cause these predictions. $\square$

# B  Fully Observed Amortized Causal Discovery

## B.1  Experimental Details

In our experiments, the amortized encoder implements two edge-propagation steps along the causal graph and consists of fully-connected networks (MLPs) or 1D CNNs with attentive pooling. The decoder implements a single edge-propagation step and consists of an MLP. Our model is optimized using ADAM [24] and implemented in PyTorch [37].

We measure our causal discovery performance by area under the receiver operator curve (AUROC), taken across all the adjacency matrices of the test set flattened and concatenated together.

### B.1.1  Datasets

**Physics Simulations**   To generate these simulations, we follow the description of the underlying physics of Kipf et al. [27] for the phase-coupled oscillators (Kuramoto) [29] and the particles connected by springs. In contrast to their simulations, however, we allow the connectivity matrix, which describes which time-series influences another, to be asymmetric. This way, it describes causal relations instead of correlations.

For both datasets, we simulate systems with $N = 5$ time-series. Our training and validation samples consist of $T = 49$ time-steps, while the test-samples are $T = 99$ time-steps long. This increased length allows us to infer causal relations on the first half of the data, and to test the future prediction performance on the second half (with $k = \{1, ..., 49\}$).

**Netsim**   The Netsim dataset simulates blood-oxygen-level-dependent (BOLD) imaging data across different regions within the human brain and is described in Smith et al. [47]. The task is to infer the directed connections, i.e. causal relations, between different brain areas.

The Netsim dataset includes simulations with different numbers of brain regions and different underlying connectivity matrices. In our experiments, we use the data from the third simulation Sim-3.mat as provided by Khanna and Tan [23]. It consists of samples from 50 subjects, each with the same underlying causal graph, each of length $T = 200$ and including $N = 15$ different brain regions. Note, that we report worse results than Khanna and Tan [23], since we assume self-connectivity for all time-series and only evaluate the causal discovery performance between *different* time-series.

The dataset is very small (50 samples) and due to this, we do not use a training/validation/test split, but use the same 50 points at each phase instead. While this is not standard machine learning practice, it still facilitates a fair comparison to the other methods, each of which are fit to individual test points. The purpose of including experiments on this dataset is not to demonstrate generalization ability, but rather to show that our method is flexible enough to work reasonably well even in the classical causal discovery setting (with one shared causal graph, and fitting the model on the test set).

### B.1.2  Architecture and Hyperparameters

We implement the MLP Encoder, CNN Encoder and MLP Decoder as described for the NRI model [27]. We use the CNN Encoder on the Kuramoto dataset and the MLP Encoder on both the particles dataset and Netsim. Their latent dimension is set to size 256.

Following our causal formulation of the NRI model, we implement Eq. (12) by masking out the values of the corresponding edges. Thus, the ordering of the edge types is not arbitrary in our setting. When conducting test-time adaptation as described in Eq. (11), we model a distribution over $\mathbb{G}$ using a non-amortized variational distribution $q(\boldsymbol{z})$ with its initial values sampled from a unit Gaussian.

We did no hyperparameter optimization for model training, but used the settings as described for the NRI model [27]. We optimize our model using ADAM [24] with a learning rate of 0.0005. In the experiments on the particles dataset, the learning rate is decayed by a factor of 0.5 every 200 epochs. We predict $k = \{1, ..., 10\}$ steps into the future. We set our batchsize to 128 and train for 500 epochs. The temperature of the Gumbel-Softmax is set to $\tau = 0.5$. During testing, this concrete distribution is replaced by a categorical distribution to obtain discrete edge predictions.

When we report the variance on the ACD-NRI results, we collected these across five different random seeds. Baselines in Kuramoto/Netsim use three seeds each, except for NGC, which uses only one due to a longer runtime (the confidence intervals shown for NGC are confidence intervals on the AUROC itself, whereas all other confidence intervals are based on variance of AUROC across seeds).

There was no thorough hyperparameter optimization done for TTA, but since there was no pre-existing implementation, some hand-tuning was performed. We use a learning rate of 0.1 for the Kuramoto and particles datasets and 0.01 for Netsim. For each, we run 1000 iterations.

### B.1.3 Baselines

We compare ACD-NRI against several baselines:

**Neural Granger Causality** From Tank et al. [49], we optimized an MLP or LSTM to do next step prediction on a sample. We found that the MLP worked best. The causal links are wherever an input weight is non-zero. We used ADAM and then line search to find exact zeros. In this method, we calculate AUROC by running with a range of sparsity hyperparameters ($\lambda = [0, 0.1, 0.2, 0.4, 0.8]$ for Kuramoto and $\lambda = [0, 0.1, 0.15, 0.2, 0.3, 0.4, 0.5, 0.6, 1]$ for Netsim). As in Tank et al. [49], we calculate a score $s$ for each edge, where $s = \min\{\lambda : z_{ij,0} = 1\}$, and use that score to calculate AUROC. Code was used from `https://github.com/icc2115/Neural-GC`.

**ESRU** Khanna and Tan [23] take a similar approach to Tank et al. [49], but they use economy statistical recurrent units (eSRU), instead of LSTMs. We found one layer worked best, and used their hyperparameters otherwise. We use sparsity hyperparameters $[0.1, 0.2, 0.3, 0.4, 0.5]$ for Kuramoto, and $[0, 0.1, 0.15, 0.2, 0.3, 0.4, 0.5, 0.6, 1]$ for Netsim. Code was used from `https://github.com/sakhanna/SRU_for_GCI`.

**MPIR** Wu et al. [52] determine where causal links exist by examining the predictive performance change when noise is added on an input variable. Code for this method and the baselines below was used from `https://github.com/tailintalent/causal`.

**Transfer Entropy** Schreiber [45] suggest this entropy-like measure between two variables to produce a metric which is likely to be higher when a causal connection exists. We use the implementation by Wu et al. [52].

**Mutual Information** Using the implementation by Wu et al. [52], we calculate the mutual information between every pair of time series.

**Linear Granger Causality** Using the implementation by Wu et al. [52], this is a linear version of Granger causality where non-zero linear weights are taken as greater causal importance.

We did not run the baselines on the particles dataset since it is two-dimensional and most baselines did not provide an obvious way for handling multi-dimensional time series. When training ACD-NRI on the particles and Kuramoto datasets, we additionally input the velocity (and phase for Kuramoto) of the time-series. Since our chosen NRI encoders and decoders are not recurrent we cannot recover this information in any other way in this model. This enables a more fair comparison to the recurrent methods, which are able to aggregate this information over several time steps.

## B.2 Additional Experimental Result - Training Curves

Fig. 7 shows the training curves when training on 100 training samples of the particles dataset. We observe that the encoder overfits on the training samples, as indicated by the AUROC performance. In contrast, the decoder shows less overfitting as indicated by the negative log-likelihood (NLL) performance.
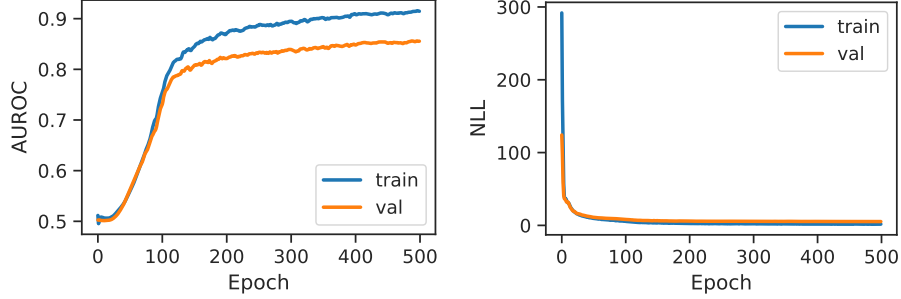


Figure 7: Training curves when training on 100 samples of the particles dataset. The encoder performance (AUROC - left) shows stronger signs of overfitting than the decoder performance (NLL - right).

## C  Amortized Causal Discovery with Unobserved Variables

### C.1  Temperature Experiments

**Implementation Details**   In this experiment, we use the CNN encoder and a simulation decoder matching the true generative ODE process. Our optimization scheme is the same as before.

For the prediction of the latent temperature, we output a uniform distribution as our posterior $q_{\phi_c}(c|\boldsymbol{x})$. One tricky aspect about this is the KL-Divergence:

$$KL(q_{\phi_c}(c|\boldsymbol{x})||p(c)) = -\int q_{\phi_c}(c|\boldsymbol{x})\log\frac{q_{\phi_c}(c|\boldsymbol{x})}{p(c)}d\boldsymbol{z} \quad . \tag{13}$$

We must ensure that our posterior support is a subset of our prior support. Otherwise, the KL-Divergence is undefined and optimization impossible. Recall that our prior is a uniform distribution over $[0, 4\alpha]$.

We output two latent parameters $a, b \in \mathbb{R}$ for each input and use these values to parametrize a mean $m$ and a half-width $w$ for the uniform distribution. First, we bound these values to represent a uniform distribution $u_1$ in $[0, 1]$. To achieve this, we let $m_1 = \sigma(a)$ and $w_1 = \sigma(b) * \min(m_1, 1 - m_1)$ with $\sigma(x) = \frac{1}{1+\exp(-x)}$. We then sample a temperature $\hat{c}_1 \sim u_1 = U(m_1 - w_1, m_1 + w_1)$, which is guaranteed to be bounded within $[0, 1]$. Stopping gradients, we use this temperature sample in the encoder $q_\phi(\boldsymbol{z}|\boldsymbol{x}, c)$ to improve the causal discovery performance.

Next, we scale this result to the desired interval $[0, 4\alpha]$. To achieve this, we feed the scaled temperature $\hat{c} = 4\alpha\hat{c}_1$ into the decoder, and use the scaled distribution $u = U(4\alpha m_1 - 4\alpha w_1, 4\alpha m_1 + 4\alpha w_1)$ to find our KL term. We allow gradients to flow through the temperature sample in both the decoder and the distribution in the KL term, which informs our parameter updates.

**Results**   Similarly to Fig. 5, we show the future prediction performance in MSE across different values of $\alpha$ in Fig. 8. Again we find a slight improvement in performance when using ACD-NRI with *Latent* compared to the baselines, although this is a noisier indicator.

Additionally, we evaluate how well the introduced latent variable learns to predict the unobserved temperature. To do so, we use the mean of the predicted posterior uniform distribution. When a discrete categorical prediction is needed for evaluation, we quantize our results into three bins based on their distance in log-space. To calculate AUROC in this three category ordinal problem, we average the AUROC between the two binary problems: category 1 vs not category 1, and category 3 vs not category 3 (category 2 vs not 2 is not a valid regression task for the purposes of AUROC
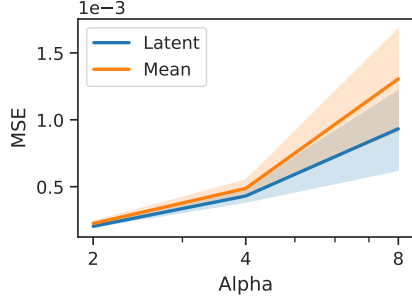
Figure 8: MSE (lower better) averaged across 5 random seeds for hidden temperature experiment. MSE for *None* baseline was much worse with MSE = 0.009, 0.02, 0.04 for $\alpha = 2, 4, 8$ (not shown in plot).
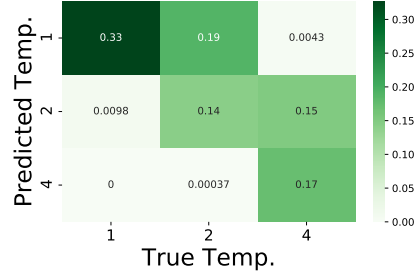
Figure 9: Confusion matrix for latent temperature prediction with $\alpha = 2$. ACD-NRI with Latent's prediction tends to be conservative: it is more likely to predict a too low temperature than a too high one.

which is concerned with ordering, since it is the middle temperature and hence the labels would not be linearly separable).

The confusion matrix between true and predicted temperature in Fig. 9 indicates that ACD-NRI with Latent's prediction tends to be conservative: it is more likely to predict a too low temperature than a too high one. This is probably due to higher temperatures incurring larger MSEs, since higher temperature systems are more chaotic and thus less predictable.

Table 3 lists the temperature prediction results across all tested values of $\alpha$. We find that we can predict the unobserved temperature quite well, especially with respect to ordering (as measured by correlation and AUROC).

|  | $\alpha$ | | |
| --- | --- | --- | --- |
|  | 2 | 4 | 8 |
| Correlation | 0.888 | 0.844 | 0.661 |
| Accuracy | 0.644 | 0.384 | 0.346 |
| AUROC (1vAll) | 0.966 | 0.935 | 0.843 |

Table 3: Latent Temperature Prediction Metrics. We treat the mean of the outputted interval of the uniform posterior as the predicted temperature. For accuracy, this value discretized in log space to get a ternary prediction. *AUC (1vAll)* averages the two one-vs-all AUC values which can be calculated in a 3-category ordinal problem.

## C.2 Unobserved Time-series

**Implementation Details** For modeling the unobserved time-series, we employ a two-layered, bi-directional long short-term memory (LSTM) [20] with a latent dimension of size 256.

**Results** The full evaluation of our experiments with an unobserved time-series can be found in Table 4. They indicate that our proposed method *ACD-NRI with latent* predicts the trajectory of the unobserved time-series (unobserved MSE) more accurately than the *Mean* imputation baseline. Even though this prediction is worse than for the *Supervised* baseline, ACD-NRI with *Latent* manages to recover the performance of the fully *Observed* baseline better than the *None* and the *Mean* imputation baselines.

Fig. 10 shows the performance of the tested methods dependent on the number of time-series that are influenced by the unobserved one. In addition to Fig. 5 in our Experiments section, these plots show the achieved accuracy and MSE results. The general trends are the same. Fig. 11 shows example trajectories and the corresponding predictions for all tested methods.

| Method | AUROC | Accuracy | MSE | unobserved MSE |
|---|---|---|---|---|
| Observed | (0.99) | (0.993) | (0.00301) | - |
| Supervised | 0.982 | 0.931 | 0.00822 | 0.0164 |
| None | 0.946 | 0.882 | 0.0119 | - |
| Mean | 0.951 | 0.881 | 0.0106 | 0.0397 |
| ACD-NRI with latent | 0.979 | 0.918 | 0.00747 | 0.0375 |

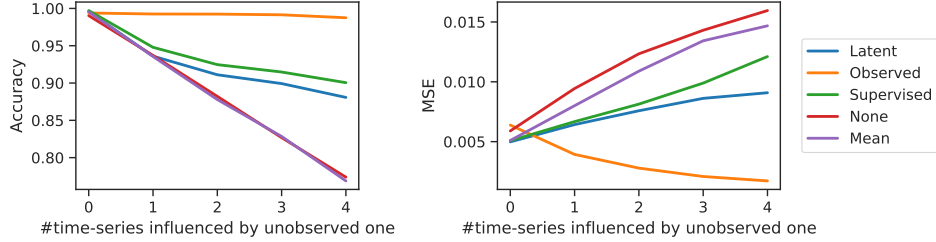Table 4: Experiments with an unobserved time-series.



Figure 10: Experiments with an unobserved particle. Performance of the various methods depends strongly on how many observed particles are influenced by the unobserved one (x-axis). The more particles that are influenced by the unobserved particle, the stronger the benefit of using an additional *Latent* variable for modeling its effects. Left - causal relation prediction accuracy (higher = better), right - MSE (lower = better).

**Uninfluenced Influencer**   Predicting the trajectory of a time-series that influences only a small number of observed time-series and is (invisibly) influenced by them is arguably very difficult. In this follow-up experiment, we reduce the difficulty of this problem by adding two assumptions: (1) the unobserved time-series influences *all* observed time-series and (2) it is not influenced by any of the observed time-series. This way, we gain more information about its trajectory (due to (1)) and its trajectory becomes easier to predict (due to (2)). Indeed, in this setup, *ACD-NRI with latent* manages to almost completely recover the performance of the fully observed baseline (Table 5). In contrast, the performance of the *None* and *Mean* imputation baselines worsens considerably in this setting. Now, all time-series are influenced by the unobserved particle - making their prediction harder when not taking into account this hidden confounder. Fig. 12 shows example trajectories and the corresponding predictions for all tested methods in this setting.

| Method | AUROC | Acuracy | MSE | unobserved MSE |
|---|---|---|---|---|
| Observed | (1.0) | (0.997) | (0.0193) | - |
| Supervised | 1.0 | 0.993 | 0.024 | 0.000615 |
| None | 0.829 | 0.76 | 0.0431 | - |
| Mean | 0.853 | 0.782 | 0.0365 | 0.0357 |
| ACD-NRI with latent | 1.0 | 0.994 | 0.0251 | 0.137 |

Table 5: Experiments with an unobserved time-series that influences all observed time-series, but is not influenced by them.
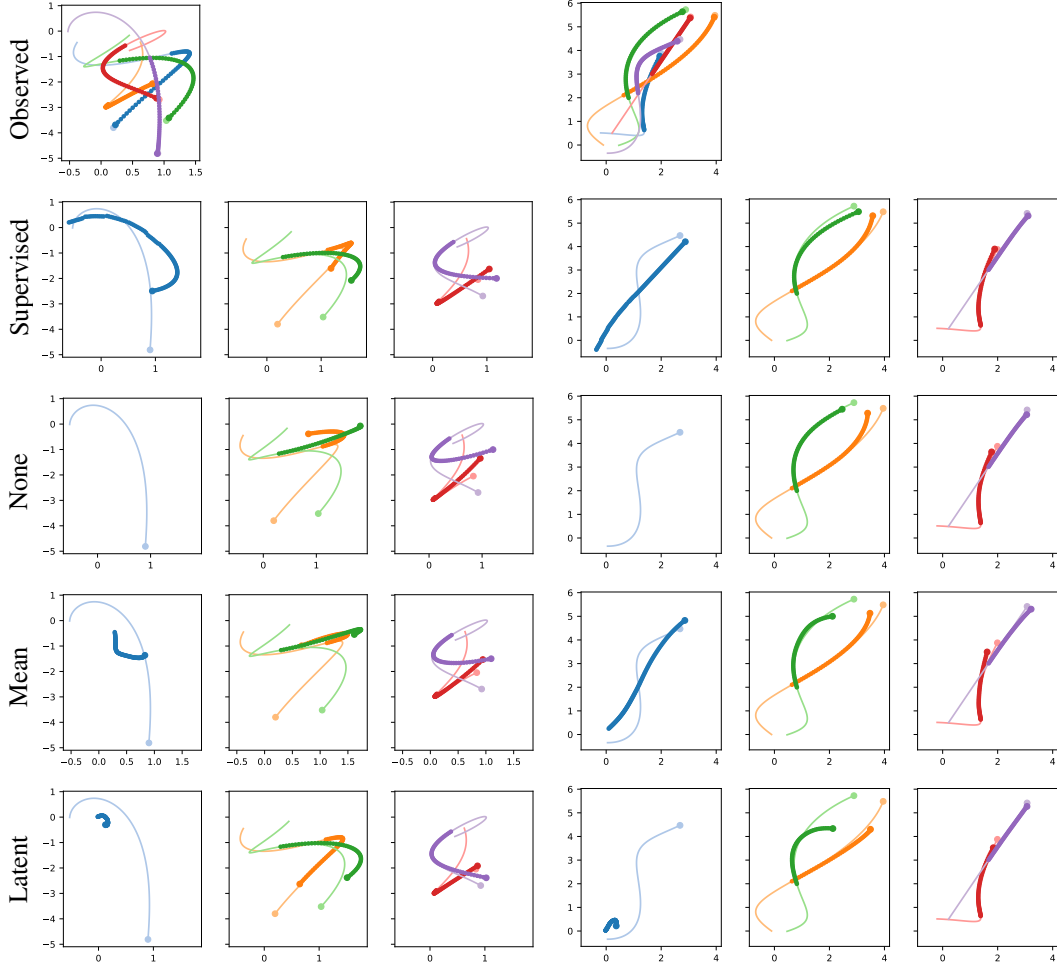
Figure 11: Predicted trajectories for all tested methods in the unobserved time-series experiment for two samples (left/right). From top to bottom: Baselines – observed, supervised, none and mean; proposed ACD-NRI with latent. The faded lines depict the ground truth trajectory; bold lines are the trajectories predicted by the model and they start after initializing the model using first half of the ground truth. Dots denote the end of the trajectories. Except for the fully observed baseline, the first panel shows the ground truth and prediction for the unobserved time-series. The second panel shows the trajectories of all time-series that are directly influenced by the unobserved one. The third panel shows the trajectories of all time-series that are not directly influenced by the unobserved one.
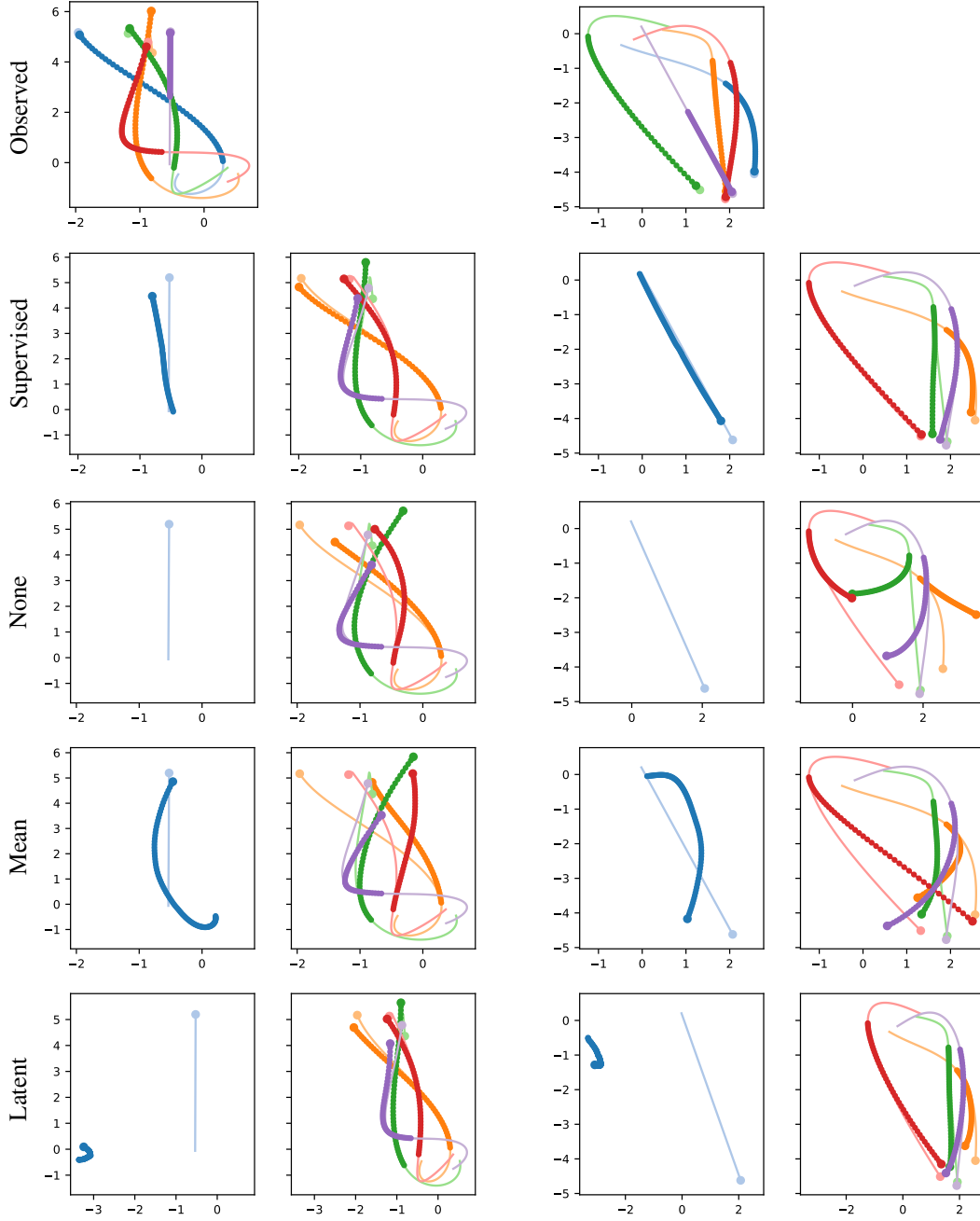
Figure 12: Predicted trajectories for all tested methods when the unobserved time-series influences all observed ones, but stay uninfluenced itself for two samples (left/right). From top to bottom: Baselines – observed, supervised, none and mean; proposed ACD-NRI with latent. The faded lines depict the ground truth trajectory; bold lines are the trajectories predicted by the model and they start after initializing the model using first half of the ground truth. Dots denote the end of the trajectories. Except for the fully observed baseline, the first panel shows the ground truth and prediction for the unobserved time-series. The second panel shows the trajectories of all observed time-series (which are all influenced by the unobserved one).