

A GRANGER CAUSALITY OF ACD

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.

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

Proof. According to Definition 2.1, time-series i does not cause j , if g_j is invariant to $\mathbf{x}_i^{\leq t}$. 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 $\mathbf{x}_i^{\leq t}$, and i does not Granger cause these predictions. \square

B CONSISTENCY

In this section, we discuss the consistency of ACD, and show that given the correct dynamics and under some assumptions, ACD learns the correct causal graph by minimizing the predictive loss ℓ . Then, we will provide some intuition on the role of the regularization term r for jointly learning the graph and dynamics. We note that the practicality of our approach will ultimately depend on the exact parameterization chosen for the encoder and decoder, as well as the forms of the loss function and regularization.

We define the domain of the observations as \mathbb{X} , letting a full, observed sequence $\mathbf{x}^{\leq t} = \mathbb{X}^{\leq t}$, dropping the s subscript for readability; the domain of the observations plus the predicted step is $\mathbb{X}^{t+1} \ni [\mathbf{x}^{\leq t}, \mathbf{x}^{t+1}]$. We recall also the domain of the graph: $\mathcal{G} \in \mathbb{G}$. Let $G \in \{0, 1\}^{n \times n}$ be the adjacency matrix describing \mathcal{G} – we assume without loss of generality that there are only two edge types (edge or no-edge), and so a binary adjacency matrix is sufficient to describe \mathcal{G} . We consider a data distribution of interest P , which is defined over $\mathbb{X}^{\leq t} \times \mathbb{G}$, and whose density function is non-zero for a set of $\mathbf{x}^{\leq t} \in \mathbb{X}^{\leq t}$ of measure greater than 0. We assume that given an observed sequence, the causal discovery problem has a correct answer — that is, $\forall \mathbf{x}^{\leq t} \in \mathbb{X}^{\leq t}$, there is at most one $\mathcal{G} \in \mathbb{G}$ such that $P(\mathbf{x}^{\leq t}, \mathcal{G}) > 0$.

As written previously, we assume that some fixed function g describes the dynamics of *all* samples $\mathbf{x} \in \mathbb{X}$ given their past observations $\mathbf{x}^{\leq t}$ and their underlying causal graph \mathcal{G} , with some additive noise ϵ^t drawn i.i.d with mean zero:

$$\mathbf{x}^{t+1} = g(\mathbf{x}^{\leq t}, \mathcal{G}) + \epsilon^t . \quad (15)$$

We suppose that we have learned the correct dynamics, i.e. our decoder f_θ is equal to g . Suppose we hope to minimize the mean squared error (MSE) of our predictions. We want to know how much additional MSE we incur by inferring the wrong graph $\hat{\mathcal{G}}$. For a sequence of length T , this MSE loss is

$$\frac{1}{T} \sum_{t=1}^T [g(\mathbf{x}^{\leq t}, \mathcal{G}) + \epsilon^t - g(\mathbf{x}^{\leq t}, \hat{\mathcal{G}})]^2 \quad (16)$$

$$= \frac{1}{T} \sum_{t=1}^T [g(\mathbf{x}^{\leq t}, \mathcal{G}) - g(\mathbf{x}^{\leq t}, \hat{\mathcal{G}})]^2 + (\epsilon^t)^2 + 2\epsilon_t [g(\mathbf{x}^{\leq t}, \mathcal{G}) - g(\mathbf{x}^{\leq t}, \hat{\mathcal{G}})] . \quad (17)$$

As $T \rightarrow \infty$, this sample mean approaches its expectation:

$$= \mathbb{E}[[g(\mathbf{x}^{\leq t}, \mathcal{G}) - g(\mathbf{x}^{\leq t}, \hat{\mathcal{G}})]^2] + \mathbb{E}[(\epsilon^t)^2] + 2\mathbb{E}[\epsilon_t] \mathbb{E}[g(\mathbf{x}^{\leq t}, \mathcal{G}) - g(\mathbf{x}^{\leq t}, \hat{\mathcal{G}})] \quad (18)$$

$$= \mathbb{E}[[g(\mathbf{x}^{\leq t}, \mathcal{G}) - g(\mathbf{x}^{\leq t}, \hat{\mathcal{G}})]^2] + \mathbb{E}[(\epsilon^t)^2] . \quad (19)$$

Since, for a long enough noisy sequence, $\mathbb{E}[(\epsilon^t)^2]$ is constant, we minimize the MSE by minimizing $\mathbb{E}[[g(\mathbf{x}^{\leq t}, \mathcal{G}) - g(\mathbf{x}^{\leq t}, \hat{\mathcal{G}})]^2]$. This term is clearly minimized at 0 for $\hat{\mathcal{G}} = \mathcal{G}$. The remaining question is, are there any $\hat{\mathcal{G}} \neq \mathcal{G}$ which also minimize this term at 0? What we will show in the next part of this

proof is that there is indeed a set of $\mathbf{x}^{\leq t}$ for which such a $\hat{\mathcal{G}} \neq \mathcal{G}$ exists, but that this set is of measure 0, and as such will occur with probability 0 under an absolutely continuous probability distribution.

To show this, we consider the noiseless dynamics:

$$\mathbf{x}^{t+1} = g(\mathbf{x}^{\leq t}, \mathcal{G}) \quad , \quad (20)$$

and suppose again that we have learned the correct dynamics, i.e. our decoder f_θ is equal to g . We let ℓ be a loss function which is minimized at zero when its two inputs are equal (such as squared error), and see that if $\hat{\mathcal{G}} = \mathcal{G}$, then $\ell(\mathbf{x}^{t+1}, \hat{\mathbf{x}}^{t+1}) = \ell(\mathbf{x}^{t+1}, g(\mathbf{x}^{\leq t}, \hat{\mathcal{G}})) = \ell(\mathbf{x}^{t+1}, g(\mathbf{x}^{\leq t}, \mathcal{G})) = \ell(\mathbf{x}^{t+1}, \mathbf{x}^{t+1}) = 0$.

What remains to show is that no graph $\hat{\mathcal{G}} \neq \mathcal{G}$ can minimize this loss function. To show this, we will show that for any two graphs $\mathcal{G}, \mathcal{G}'$, for which $\hat{\mathcal{G}} \neq \mathcal{G}$, the subset of $\mathbf{x}^{\leq t}$ for which $g(\mathbf{x}^{\leq t}, \mathcal{G}) = g(\mathbf{x}^{\leq t}, \mathcal{G}')$ has measure zero. Therefore, the probability of sampling $\mathbf{x}^{\leq t}$ for which $g(\mathbf{x}^{\leq t}, \mathcal{G}) = g(\mathbf{x}^{\leq t}, \mathcal{G}')$ is zero. We then note that at most 3^n of these sets exist, yielding a total measure of 0.

First, we assume that both the input and output of g are in continuous spaces, and we make two assumptions on g :

1. If $G_{ij} = 0$, then $\forall \mathbf{x}^{\leq t}, \tilde{\mathbf{x}}^{\leq t} \in \mathbb{X}^{\leq t}$ which are equal except for the i -th time-series, $g_j(\mathbf{x}^{\leq t}, \mathcal{G}) = g_j(\tilde{\mathbf{x}}^{\leq t}, \mathcal{G})$.
2. Let $\mathbf{x}^{\leq t} \in \mathbb{X}^{\leq t}, \mathcal{G} \in \mathcal{G}$. Suppose $G_{ij} = 1$. Define the function $g_j^i(v; k) = g(\mathbf{x}_{v \rightarrow \{i\}}, \mathcal{G})$, where $\mathbf{x}_{v \rightarrow \{i\}}$ refers to a perturbation of the i -th time-series in $\mathbf{x}^{\leq t}$ such that \mathbf{x}_i^{t-k} takes on the value v , with every other position remaining the same. Then, for some k , g_j^i is strictly monotonic for all i, j . Going forward, we use the k for which this holds and simply let $g_j^i(v) = g_j^i(v; k)$.

The first assumption follows directly from the definition of non-linear Granger causality (Definition 2.1), and describes the same phenomenon as proven for our decoder in Appendix A: when there is no edge in the causal graph between time-series i and j , g_j is invariant to $\mathbf{x}_i^{\leq t}$. The second assumption describes that, when there *is* an edge in the causal graph, the j -th output of g_j will be sensitive to changes in i in a predictable way.

Suppose we have some $P(\mathbf{x}^{\leq t}, \mathcal{G}) > 0$. Let j be an output dimension, and let $\mathcal{I}, \mathcal{I}' \subseteq \{1 \dots n\}$ be sets of indices with $\mathcal{I} \cap \mathcal{I}' = \emptyset$ such that $G_{ij} = 0 \forall i \in \mathcal{I}$ and $G_{ij} = 1 \forall i \in \mathcal{I}'$. Let $\mathbf{x}_{\mathcal{I}}^{\leq t}, \mathbf{x}_{\mathcal{I}'}^{\leq t}$ be the subsets of positions in $\mathbf{x}^{\leq t}$ indexed by the corresponding index set, and $\mathbf{x}_{\mathcal{I},k}^{\leq t}, \mathbf{x}_{\mathcal{I}',k}^{\leq t}$ the k -th position in each.

As in the second assumption above, the notation $\mathbf{x}_{v \rightarrow \mathcal{I}, w \rightarrow \mathcal{I}'}^{\leq t}$ will refer to the assignment of values v to indices \mathcal{I} and the assignment of values w to indices \mathcal{I}' , both at the k -th position (v and w may be $|\mathcal{I}|$ and $|\mathcal{I}'|$ -length vectors for this reason). Let $g_j^{\rightarrow}(v, w, \mathcal{G}) = g_j(\mathbf{x}_{v \rightarrow \mathcal{I}, w \rightarrow \mathcal{I}'}^{\leq t}, \mathcal{G})$. We will let $\mathbb{X}_{\mathcal{I},k}$ be the domain of vectors containing values for all k -th positions of time-series indexed by \mathcal{I} , and $\mathbb{X}_{\mathcal{I}',k}$ analogously.

We begin by noting that we can set the k -th position values in \mathcal{I} to any values v without changing the j -th output of g , following from the first assumption on g . Formally, $\forall v \in \mathbb{X}_{\mathcal{I},k}, g_j(\mathbf{x}^{\leq t}, \mathcal{G}) = g_j^{\rightarrow}(v, \mathbf{x}_{\mathcal{I}',k}^{\leq t}, \mathcal{G}) = c$ for some $c \in \mathbb{R}$.

Now, we consider the graph $\mathcal{G}^{\mathcal{I}j=1}$, which is defined as $G_{il}^{\mathcal{I}j=1} = 1$ if $i \in \mathcal{I}, l = j$, and G_{il} otherwise. Intuitively, this is just the adjacency matrix resulting from adding edges at all of the indices in \mathcal{I} to \mathcal{G} . We note that $g_j^{\rightarrow}(v, \mathbf{x}_{\mathcal{I}',k}^{\leq t}, \mathcal{G}^{\mathcal{I}j=1})$ is monotonic in v by the second assumption on g . In particular, this means that the set of $\mathbf{x}^{\leq t}$ for which $g_j^{\rightarrow}(\mathbf{x}_{\mathcal{I},k}^{\leq t}, \mathbf{x}_{\mathcal{I}',k}^{\leq t}, \mathcal{G}^{\mathcal{I}j=1}) = g_j(\mathbf{x}^{\leq t}, \mathcal{G}^{\mathcal{I}j=1}) = c = g_j(\mathbf{x}^{\leq t}, \mathcal{G})$ is small; it is a curve of dimension $|\mathcal{I}| - 1 < n$ and as such has measure 0. We can see this because this set of \mathbf{x} is the intersection of a monotonic curve with the constant c , where the curve is $\{(y, v) | g_j^{\rightarrow}(v, \mathbf{x}_{\mathcal{I}',k}^{\leq t}, \mathcal{G}^{\mathcal{I}j=1}) = y\}$, which is $|\mathcal{I}|$ -dimensional since that is the dimensionality of v . We can use similar logic to address the possibility of flipping 1's to 0's in the adjacency matrix; this will similarly change the prediction with probability 1.

Now that we have flipped all the edges in \mathcal{I} from 0 to 1, we turn to additionally flipping all the edges in \mathcal{I}' from 1 to 0. We consider the graph $\mathcal{G}^{\mathcal{I}j=1, \mathcal{I}'j=0}$, whose adjacency matrix is

defined as $G_{il}^{\mathcal{I}j=1, \mathcal{I}'j=0} = 1$ if $i \in \mathcal{I}, l = j$, $G_{il}^{\mathcal{I}j=1, \mathcal{I}'j=0} = 0$ if $i \in \mathcal{I}', l = j$, and G_{il} otherwise. By the first assumption on g above, we have that $\forall w \in \mathbb{X}_{\mathcal{I}', k}, g_j(\mathbf{x}^{\leq t}, \mathcal{G}^{\mathcal{I}j=1, \mathcal{I}'j=0}) = g_j^{\rightarrow}(\mathbf{x}_{\mathcal{I}, k}^{\leq t}, w, \mathcal{G}^{\mathcal{I}j=1, \mathcal{I}'j=0}) = c'$ for some $c' \in \mathbb{R}$.

Now we ask, what is the set of $\mathbf{x}^{\leq t}$ where $g_j(\mathbf{x}^{\leq t}, \mathcal{G}^{\mathcal{I}j=1, \mathcal{I}'j=0}) = g_j(\mathbf{x}^{\leq t}, \mathcal{G})$? We note that, in general, for all v, w , we have

$$g_j^{\rightarrow}(\mathbf{x}_{\mathcal{I}, k}^{\leq t}, w, \mathcal{G}^{\mathcal{I}j=1, \mathcal{I}'j=0}) = g_j(\mathbf{x}^{\leq t}, \mathcal{G}^{\mathcal{I}j=1, \mathcal{I}'j=0}) \stackrel{?}{\neq} g_j(\mathbf{x}^{\leq t}, \mathcal{G}) = g_j^{\rightarrow}(v, \mathbf{x}_{\mathcal{I}', k}^{\leq t}, \mathcal{G}) \quad (21)$$

We want to know, for how many values of v, w do we have

$$g_j^{\rightarrow}(v, w, \mathcal{G}^{\mathcal{I}j=1, \mathcal{I}'j=0}) = g_j^{\rightarrow}(v, w, \mathcal{G}) \quad (22)$$

This is the intersection of two curves — the first is $|\mathcal{I}'|$ -dimensional and the second is $|\mathcal{I}|$ -dimensional. Therefore, their intersection is at most $\min(|\mathcal{I}|, |\mathcal{I}'|)$ -dimensional, which is $\leq \frac{n}{2}$. This has also measure 0. \square

The regularization term r is important for tackling the remaining challenge of jointly learning both the dynamics and the causal graph. Without it, the loss cannot differentiate between a model that represents a non-causal relation correctly in $\hat{\mathcal{G}}$, or incorrectly as an empty edge-type within f_{θ} . By regularizing $\hat{\mathcal{G}}$, this symmetry can be broken and the loss can enforce the correct solution. Additionally, the regularization term allows us to include prior knowledge into our model about what we think the true sparsity of the graph should be.

C FULLY OBSERVED AMORTIZED CAUSAL DISCOVERY

C.1 EXPERIMENTAL DETAILS

C.1.1 DATASETS

Physics Simulations To generate these simulations, we follow the description of the underlying physics of Kipf et al. (2018) for the phase-coupled oscillators (Kuramoto) (Kuramoto, 1975) 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 generate 50,000 training and 10,000 validation samples. We restrict the number of test samples to 200, since the previous methods we compare to must be refit for each individual sample. 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 $k = \{1, \dots, 49\}$).

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. (2011). 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 (2020). 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 (2020), 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).

C.1.2 ARCHITECTURE AND HYPERPARAMETERS

Our model is implemented in PyTorch (Paszke et al., 2019). We did no hyperparameter optimization for model training, but used the settings as described for the NRI model (Kipf et al., 2018). The latent dimension throughout the model is set to size 256. We optimize our model using ADAM (Kingma and Ba, 2015) 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 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.

There was no thorough hyperparameter optimization done for test-time adaptation (TTA). 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.

Encoder In our experiments, the amortized encoder applies a graph neural network $f_{enc,\phi}$ on the input. It implements two edge-propagation steps along the causal graph:

$$\psi_j^1 = f_{emb}(x_j) \quad (23)$$

$$\psi_{ij}^1 = f_e^1([\psi_i^1, \psi_j^1]) \quad (24)$$

$$\psi_j^2 = f_v^2\left(\sum_{i \neq j} \psi_{ij}^1\right) \quad (25)$$

$$\psi_{ij} = f_e^2([\psi_i^2, \psi_j^2]) \quad (26)$$

f_e^1, f_e^2 and f_v^2 are fully-connected networks (MLPs). On both the particles dataset and Netsim, f_{emb} is an MLP as well (MLP Encoder); on the Kuramoto dataset, we use a 1D CNN with attentive pooling (Lin et al., 2017) instead (CNN Encoder).

When conducting test-time adaptation as described in Eq. (7), we remove the encoder and model a distribution over \mathbb{G} using a non-amortized variational distribution $q(z)$ with its initial values sampled from a unit Gaussian.

Decoder The decoder implements a single edge-propagation step according to equations 12-14. It uses MLPs for both f_e and f_v . To improve performance, we train the decoder to predict several time-steps into the future. For this, we replace the true input x^t with the predicted μ^t for $k = 10$ steps.

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.

Since our physics simulations are differentiable, we can replace the decoder with the ground-truth dynamics and backpropagate through them. We call this setup the simulation decoder.

Variance When we report the variance on the ACD 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).

C.1.3 BASELINES

We compare ACD against several baselines:

Neural Granger Causality From Tank et al. (2018), 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. (2018), 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 (2020) take a similar approach to Tank et al. (2018), 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. (2020) 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 (2000) 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. (2020).

Mutual Information Using the implementation by Wu et al. (2020), we calculate the mutual information between every pair of time series.

Linear Granger Causality Using the implementation by Wu et al. (2020), 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 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.

C.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.

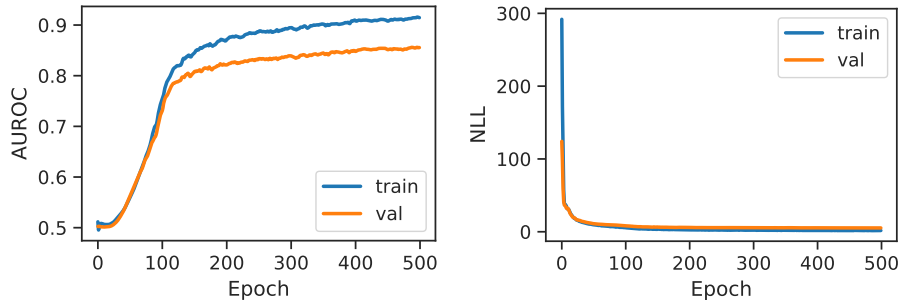


Figure 7: Training curves when training on 100 samples of the particles dataset. The encoder performance (AUROC - left) shows stronger signs of overfitting than the decoder performance (NLL - right).

D AMORTIZED CAUSAL DISCOVERY WITH UNOBSERVED VARIABLES

D.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.

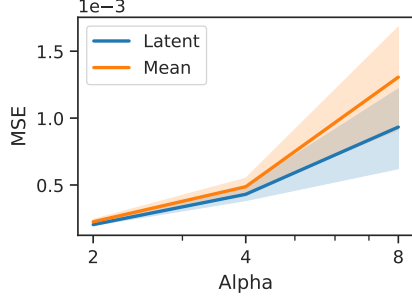


Figure 8: MSE (lower better) averaged across 5 random seeds for hidden temperature experiment. MSE for *None* baseline was much worse with MSE = 0.009, 0.02, 0.04 for $\alpha = 2, 4, 8$ (not shown in plot).

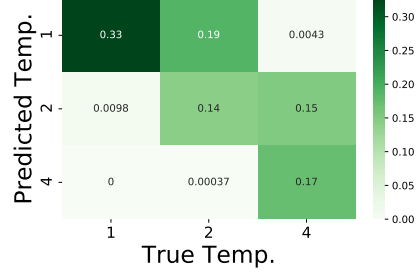


Figure 9: Confusion matrix for latent temperature prediction with $\alpha = 2$. ACD with Latent’s prediction tends to be conservative: it is more likely to predict a too low temperature than a too high one.

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

$$KL(q_{\phi_c}(c|\mathbf{x})||p(c)) = - \int q_{\phi_c}(c|\mathbf{x}) \log \frac{q_{\phi_c}(c|\mathbf{x})}{p(c)} dz . \quad (27)$$

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 $q_{\phi}(z|\mathbf{x}, c)$ 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.

Additional Results Similarly to Fig. 5, we show the future prediction performance in MSE across different values of α in Fig. 8. Again we find a slight improvement in performance when using ACD 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 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 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 α . We find that we can predict the unobserved temperature quite well, especially with respect to ordering (as measured by correlation and AUROC).

	α		
	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.

D.2 UNOBSERVED TIME-SERIES

Implementation Details For modeling the unobserved time-series, we employ a two-layered, bi-directional long short-term memory (LSTM) (Hochreiter and Schmidhuber, 1997) with a latent dimension of size 256.

Additional Results The full evaluation of our experiments with an unobserved time-series can be found in Table 4. Our results indicate that our proposed method *ACD 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 with Latent* manages to recover the performance of the fully *Observed* baseline better than the *None* and the *Mean* imputation baselines.

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 with latent	0.979	0.918	0.00747	0.0375

Table 4: Experiments with an unobserved time-series.

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.

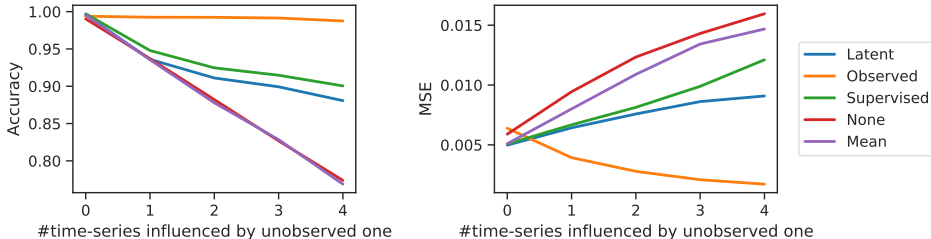


Figure 10: Experiments with an unobserved particle. Performance of the various methods depends strongly on how many observed particles are influenced by the unobserved one (x-axis). The more particles that are influenced by the unobserved particle, the stronger the benefit of using an additional *Latent*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

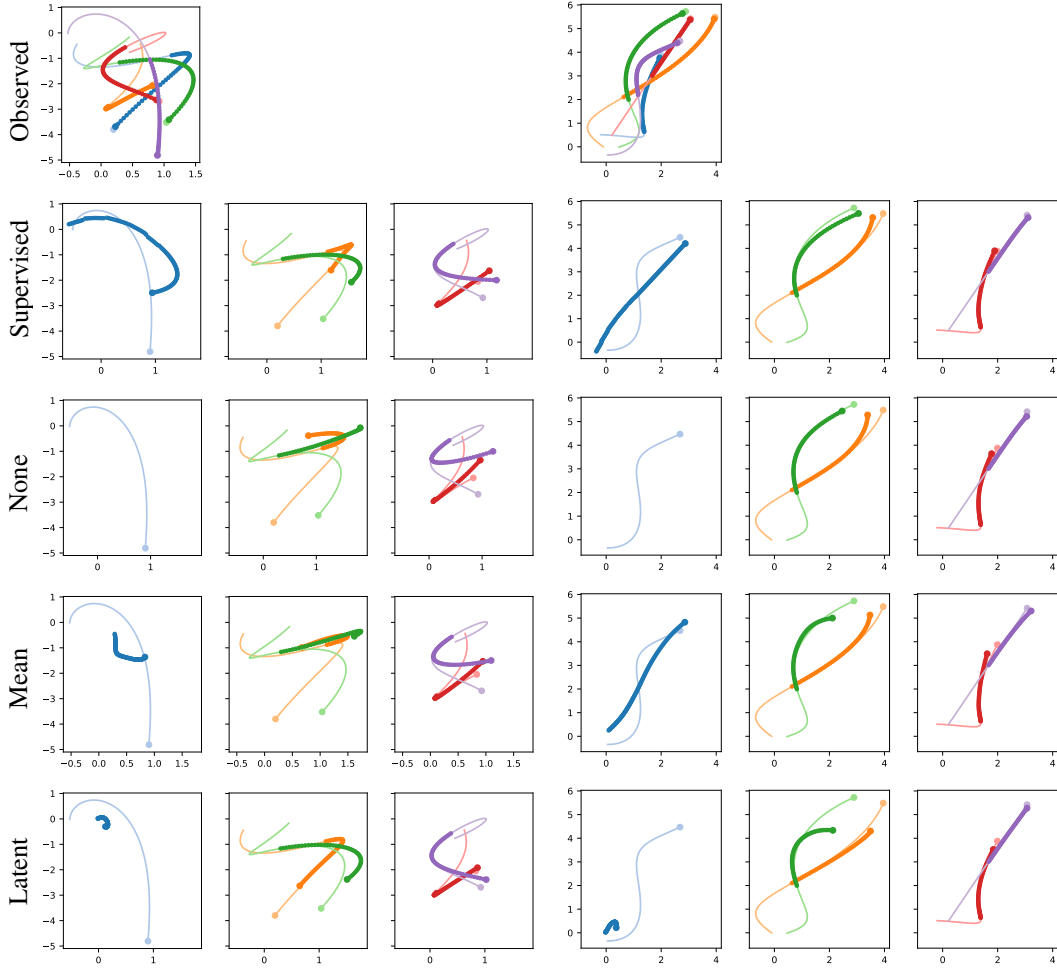


Figure 11: Predicted trajectories for all tested methods in the unobserved time-series experiment for two samples (left/right). From top to bottom: Baselines – observed, supervised, none and mean; proposed ACD with latent. The faded lines depict the ground truth trajectory; bold lines are the trajectories 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.

Additional Experiment: 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 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	Accuracy	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 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.

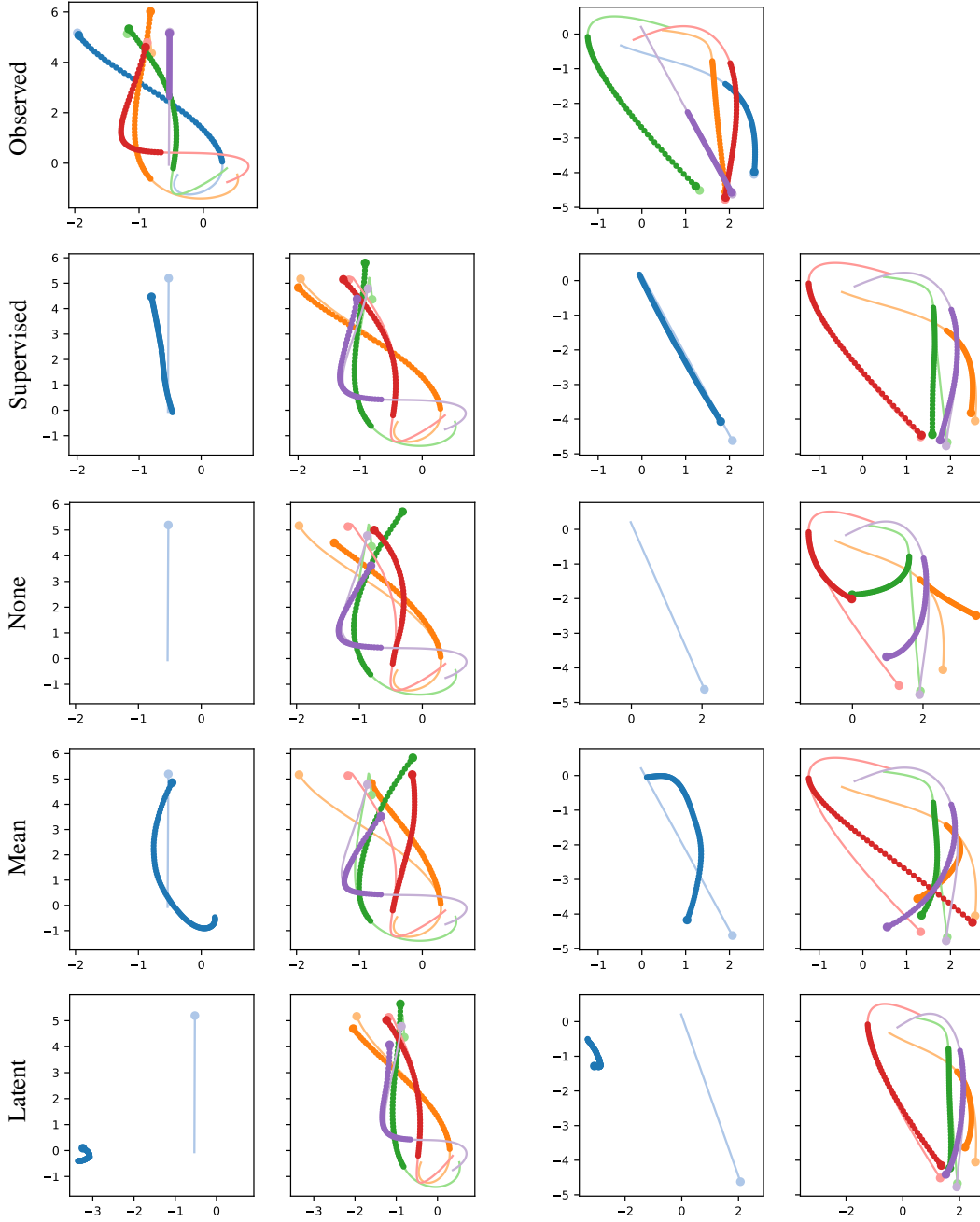


Figure 12: Predicted trajectories for all tested methods when the unobserved time-series influences all observed ones, but stay uninfluenced itself for two samples (left/right). From top to bottom: Baselines – observed, supervised, none and mean; proposed ACD with latent. The faded lines depict the ground 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).