## **Neural Granger Causality**

Alex Tank\*, Ian Covert\*, Nick Foti, Ali Shojaie, Emily B. Fox

Abstract—While most classical approaches to Granger causality detection assume linear dynamics, many interactions in real-world applications, like neuroscience and genomics, are inherently nonlinear. In these cases, using linear models may lead to inconsistent estimation of Granger causal interactions. We propose a class of nonlinear methods by applying structured multilayer perceptrons (MLPs) or recurrent neural networks (RNNs) combined with sparsity-inducing penalties on the weights. By encouraging specific sets of weights to be zero—in particular, through the use of convex group-lasso penalties—we can extract the Granger causal structure. To further contrast with traditional approaches, our framework naturally enables us to efficiently capture long-range dependencies between series either via our RNNs or through an automatic lag selection in the MLP. We show that our neural Granger causality methods outperform state-of-the-art nonlinear Granger causality methods on the DREAM3 challenge data. This data consists of nonlinear gene expression and regulation time courses with only a limited number of time points. The successes we show in this challenging dataset provide a powerful example of how deep learning can be useful in cases that go beyond prediction on large datasets. We likewise illustrate our methods in detecting nonlinear interactions in a human motion capture dataset.

Index Terms—time series, Granger causality, neural networks, structured sparsity, interpretability



#### 1 Introduction

In many scientific applications of multivariate time series, it is important to go beyond prediction and forecasting and instead interpret the structure within time series. Typically, this structure provides information about the contemporaneous and lagged relationships within and between individual series and how these series interact. For example, in neuroscience it is important to determine how brain activation spreads through brain regions [1], [2], [3], [4]; in finance it is important to determine groups of stocks with low covariance to design low risk portfolios [5]; and, in biology, it is of great interest to infer gene regulatory networks from time series of gene expression levels [6], [7]. However, for a given statistical model or methodology, there is often a tradeoff between the interpretability of these structural relationships and expressivity of the model dynamics.

Among the many choices for understanding relationships between series, Granger causality [8], [9] is a commonly used framework for time series structure discovery that quantifies the extent to which the past of one time series aids in predicting the future evolution of another time series. When an entire system of time series is studied, networks of Granger causal interactions may be uncovered [10]. This is in contrast to other types of structure discovery, like coherence [11] or lagged correlation [11], which analyze strictly bivariate covariance relationships. That is, Granger causality metrics depend on the activity of the entire system of time series under study, making them more appropriate for understanding high-dimensional complex data streams. Methodology for estimating Granger causality may be separated into two classes, model-based and model-free.

- \* Denotes equal contribution.
- Alex Tank was with the Department of Statistics, University of Washington, Seattle, WA, 98103. E-mail: alextank@uw.edu
- Ian Covert, Nicholas Foti, and Emily Fox were with the Department of Computer Science, University of Washington, Seattle, WA, 98103.
- Ali Shojaie was with the Department of Biostatistics, University of Washington, Seattle, WA, 98103

Most classical model-based methods assume linear time series dynamics and use the popular vector autoregressive (VAR) model [7], [9]. In this case, the time lags of a series have a linear effect on the future of each other series, and the magnitude of the linear coefficients quantifies the Granger causal effect. Sparsity-inducing regularizers, like the Lasso [12] or group lasso [13], help scale linear Granger causality estimation in VAR models to the high-dimensional setting [7], [14].

In classical linear VAR methods, one must explicitly specify the maximum time lag to consider when assessing Granger causality. If the specified lag is too short, Granger causal connections occurring at longer time lags between series will be missed while overfitting may occur if the lag is too large. Lag selection penalties, like the hierarchical lasso [15] and truncating penalties [16], have been used to automatically select the relevant lags while protecting against overfitting. Furthermore, these penalties lead to a sparse network of Granger causal interactions, where only a few Granger causal connections exist for each series—a crucial property for scaling Granger causal estimation to the high-dimensional setting, where the number of time series and number of potentially relevant time lags all scale with the number of observations [17].

Model-based methods may fail in real world cases when the relationships between the past of one series and future of another falls outside of the model class [18], [19], [20]. This typically occurs when there are nonlinear dependencies between the past of one series and the future. Model-free methods, like transfer entropy [2] or directed information [21], are able to detect these nonlinear dependencies between past and future with minimal assumptions about the predictive relationships. However, these estimators have high variance and require large amounts of data for reliable estimation. These approaches also suffer from a curse of dimensionality [22] when the number of series grows, making them inappropriate in the high-dimensional setting.

Neural networks are capable of representing complex,

nonlinear, and non-additive interactions between inputs and outputs. Indeed, their time series variants, such as autoregressive multilayer perceptrons (MLPs) [23], [24], [25] and recurrent neural networks (RNNs) like long-short term memory networks (LSTMs) [26] have shown impressive performance in forecasting multivariate time series given their past [27], [28], [29]. While these methods have shown impressive predictive performance, they are essentially black box methods and provide little interpretability of the multivariate structural relationships in the series. A second drawback is that jointly modeling a large number of series leads to many network parameters. As a result, these methods require much more data to fit reliably and tend to perform poorly in high-dimensional settings.

We present a framework for structure learning in MLPs and RNNs that leads to interpretable nonlinear Granger causality discovery. The proposed framework harnesses the impressive flexibility and representational power of neural networks. It also sidesteps the black-box nature of many network architectures by introducing component-wise architectures that disentangle the effects of lagged inputs on individual output series. For interpretability and an ability to handle limited data in the high-dimensional setting, we place sparsity-inducing penalties on particular groupings of the weights that relate the histories of individual series to the output series of interest. We term these sparse component-wise models, e.g. cMLP and cLSTM, when applied to the MLP and LSTM, respectively. In particular, we select for Granger causality by adding group sparsity penalties [13] on the outgoing weights of the inputs.

As in linear methods, appropriate lag selection is crucial for Granger causality selection in nonlinear approaches especially in highly parametrized models like neural networks. For the MLP, we introduce two more structured group penalties [15], [30] [31] that automatically detect both nonlinear Granger causality and also the lags of each inferred interaction. Our proposed cLSTM model, on the other hand, sidesteps the lag selection problem entirely because the recurrent architecture efficiently models long range dependencies [26]. When the true network of nonlinear interactions is sparse, both the cMLP and cLSTM approaches will select a subset of the time series that Grangercause the output series, no matter the lag of interaction. To our knowledge, these approaches represent the first set of nonlinear Granger causality methods applicable in high dimensions without requiring precise lag specification.

We first validate our approach and the associated penalties via simulations on both linear VAR and nonlinear Lorenz-96 data [32], showing that our nonparametric approach accurately selects the Granger causality graph in both linear and nonlinear settings. Second, we compare our cMLP and cLSTM models with existing Granger causality approaches [33], [34] on the difficult DREAM3 gene regulatory network recovery benchmark datasets [35] and find that our methods outperform a wide set of competitors across all five datasets. Finally, we use our cLSTM method to explore Granger causal interactions between body parts during natural motion with a highly nonlinear and complex dataset of human motion capture [36], [37]. Our implementation is available online: https://github.com/iancovert/Neural-GC.

Traditionally, the success stories of neural networks have been on prediction tasks in large datasets. In contrast, here, our performance metrics relate to our ability to produce interpretable structures of interaction amongst the observed time series. Furthermore, these successes are achieved in limited data scenarios. Our ability to produce interpretable structures and train neural network models with limited data can be attributed to our use of structured sparsity-inducing penalties and the regularization such penalties provide, respectively. We note that sparsity inducing penalties have been used for architecture selection in neural networks [38], [39]. However, the focus of the architecture selection was on improving predictive performance rather than on returning interpretable structures of interaction among observed quantities.

More generally, our proposed formulation shows how structured penalties common in regression [30], [31] may be generalized for structured sparsity and regularization in neural networks. This opens up new opportunities to use these tools in other neural network context, especially as applied to structure learning problems. In concurrent work, a similar notion of sparse-input neural networks were developed for high-dimensional regression and classification tasks for independent data [40].

#### 2 LINEAR GRANGER CAUSALITY

Let  $\mathbf{x}_t \in \mathbb{R}^p$  be a p-dimensional stationary time series and assume we have observed the process at T time points,  $(\mathbf{x}_1, \dots, \mathbf{x}_T)$ . Using a model-based approach, as is our focus, Granger causality in time series analysis is typically studied using the vector autoregressive model (VAR) [9]. In this model, the time series at time t,  $\mathbf{x}_t$ , is assumed to be a linear combination of the past K lags of the series

$$\mathbf{x}_{t} = \sum_{k=1}^{K} A^{(k)} \mathbf{x}_{t-k} + e_{t}, \tag{1}$$

where  $A^{(k)}$  is a  $p \times p$  matrix that specifies how lag k affects the future evolution of the series and  $e_t$  is zero mean noise. In this model, time series j does not Granger-cause time series i if and only if for all k,  $A^{(k)}_{ij} = 0$ . A Granger causal analysis in a VAR model thus reduces to determining which values in  $A^{(k)}$  are zero over all lags. In higher dimensional settings, this may be determined by solving a group lasso regression problem [41]

$$\min_{A^{(1)},\dots,A^{(K)}} \sum_{t=K}^{T} \|\mathbf{x}_{t} - \sum_{k=1}^{K} A^{(k)} \mathbf{x}_{t-k} \|_{2}^{2} + \lambda \sum_{ij} \|(A_{ij}^{(1)},\dots,A_{ij}^{(K)})\|_{2}, \tag{2}$$

where  $\|\cdot\|_2$  denotes the  $L_2$  norm. The group lasso penalty over all lags of each (i,j) entry,  $\|(A_{ij}^{(1)},\ldots,A_{ij}^{(K)}\|_2$  jointly shrinks all  $A_{ij}^k$  parameters to zero across all lags k [13]. The hyper-parameter  $\lambda > 0$  controls the level of group sparsity.

The group penalty in Equation (2) may be replaced with a structured hierarchical penalty [30], [42] that automatically

selects the lag of each Granger causal interaction [15]. Specifically, the hierarchical lag selection problem is given by

$$\min_{A^{(1)},\dots,A^{(K)}} \sum_{t=K}^{T} \|\mathbf{x}_{t} - \sum_{k=1}^{K} A^{(k)} \mathbf{x}_{t-k}\|_{2}^{2} + \lambda \sum_{ij} \sum_{k=1}^{K} \|(A_{ij}^{(k)}, \dots, A_{ij}^{(K)})\|_{2}, \quad (3)$$

where  $\lambda>0$  now controls the lag order selected for each interaction. Specifically, at higher values of  $\lambda$  there exists a k for each (i,j) pair such that the entire contiguous set of lags  $(A_{ij}^{(k)},\ldots,A_{ij}^{(K)})$  is shrunk to zero. If k=1 for a particular (i,j) pair, then all lags are equal to zero and series i does not Granger-cause series j; thus, this penalty simultaneously selects for Granger non-causality and the lag of each Granger causal pair.

### 3 Models for Neural Granger Causality

### 3.1 Adapting Neural Networks for Granger Causality

A *nonlinear* autoregressive model (NAR) allows  $\mathbf{x}_t$  to evolve according to more general nonlinear dynamics [25]

$$\mathbf{x_t} = g(x_{< t1}, \dots, x_{< tp}) + e_t \tag{4}$$

where  $x_{< ti} = (\dots, x_{(t-2)i}, x_{(t-1)i})$  denotes the past of series i and we assume additive zero mean noise  $e_t$ .

In a forecasting setting, it is common to jointly model the full nonlinear functions g using neural networks. Neural networks have a long history in NAR forecasting, using both traditional architectures [25], [43], [44] and more recent deep learning techniques [27], [29], [45]. These approaches either utilize an MLP where the inputs are  $x_{<t} = x_{(t-1):(t-K)}$ , for some lag K, or a recurrent network, like an LSTM.

There are two problems with applying the standard neural network NAR model in the context of inferring Granger causality. The first is that these models act as black boxes that are difficult to interpret. Due to sharing of hidden layers, it is difficult to specify sufficient conditions on the weights that simultaneously allows series j to Granger cause series i but not Granger cause series i' for  $i \neq i'$ . Second, a joint network over all  $x_{ti}$  for all i assumes that each time series depends on the same past lags of the other series. However, in practice, each  $x_{ti}$  may depend on different past lags of the other series.

To tackle these challenges, we propose a structured neural network approach to modeling and estimation. First, instead of modeling g jointly across all outputs  $x_t$ , as is standard in multivariate forecasting, we instead focus on each output *component* with a separate model:

$$x_{ti} = g_i(x_{< t1}, \dots, x_{< tp}) + e_{ti}$$

Here,  $g_i$  is a function that specifies how the past K lags are mapped to series i. In this context, Granger non-causality between two series j and i means that the function  $g_i$  does not depend on  $x_{< tj}$ , the past lags of series j. More formally, **Definition 1.** Time series j is Granger non-causal for time series i if for all  $(x_{< t1}, \ldots, x_{< tp})$  and all  $x'_{< tj} \neq x_{< tj}$ ,

$$g_i(x_{< t1}, \dots, x_{< tj}, \dots, x_{< tp}) = g_i(x_{< t1}, \dots, x'_{< tj}, \dots x_{< tp})$$

that is,  $g_i$  is invariant to  $x_{< tj}$ .

In Section 3.2 and 3.3 we consider these component-wise models in the context of MLPs and LSTMs. We examine a set of sparsity inducing penalties as in Equations (2) and (3) that allow us to infer the invariances of Definition 1 that lead us to identify Granger non-causality.

### 3.2 Sparse Input MLPs

Our first approach is to model each output component  $g_i$  with a separate MLP, so that we can easily disentangle the effects from inputs to outputs. We refer to this approach, displayed pictorially in Figure 1, as a component-wise MLP (cMLP). Let  $g_i$  take the form of an MLP with L-1 layers and let the vector  $h_t^l \in \mathbb{R}^H$  denote the values of the m-dimensional lth hidden layer at time t. The parameters of the neural network are given by weights  $\mathbf{W}$  and biases  $\mathbf{b}$  at each layer,  $\mathbf{W} = \{W^1, \dots, W^L\}$  and  $\mathbf{b} = \{b^1, \dots, b^L\}$ . To draw an analogy with the time series VAR model, we further decompose the weights at the first layer across time lags,  $W^1 = \{W^{11}, \dots, W^{1K}\}$ . The dimensions of the parameters are given by  $W^1 \in \mathbb{R}^{H \times pK}$ ,  $W^l \in \mathbb{R}^{H \times H}$  for 1 < l < L,  $W^L \in \mathbb{R}^H$ ,  $b^l \in \mathbb{R}^H$  for l < L and  $b^L \in \mathbb{R}$ . Using this notation, the vector of first layer hidden values at time t are given by

$$h_t^1 = \sigma \left( \sum_{k=1}^K W^{1k} \mathbf{x}_{t-k} + b^1 \right),$$
 (5)

where  $\sigma$  is an activation function. Typical activation functions are either logistic or tanh functions. The vector of hidden units in subsequent layers is given by a similar form, also with  $\sigma$  activation functions:

$$h_t^l = \sigma \left( W^l h_t^{l-1} + b^l \right). \tag{6}$$

After passing through the L-1 hidden layers, the time series output,  $x_{ti}$ , is given by a linear combination of the units in the final hidden layer

$$x_{ti} = g_i(x_{< t}) + e_{ti} = W^L h_t^{L-1} + b^L + e_{ti}$$
 (7)

where  $W^L$  is the linear output decoder and  $h^L_t$  is the final hidden output from the final L-1th layer. The error term,  $e_{ti}$ , is modeled as mean zero Gaussian noise. We chose this linear output decoder since our primary motivation involves real-valued multivariate time series. However, other decoders like a logistic, softmax, or poisson likelihood with exponential link function [46], could be used to model nonlinear Granger causality in multivariate binary [47], categorical [48], or positive count time series [47].

#### 3.2.1 Penalized Selection of Granger Causality in the cMLP

In Equation (5), if the jth column of the first layer weight matrix,  $W^{1k}_{:j}$ , contains zeros for all k, then series j does not Granger-cause series i. That is,  $x_{(t-k)j}$  for all k does not influence the hidden unit  $h^1_t$  and thus the output  $x_{ti}$ . Following Definition 1, we see  $g_i$  is invariant to  $x_{< tj}$ . Thus, analogously to the VAR case, one may select for Granger

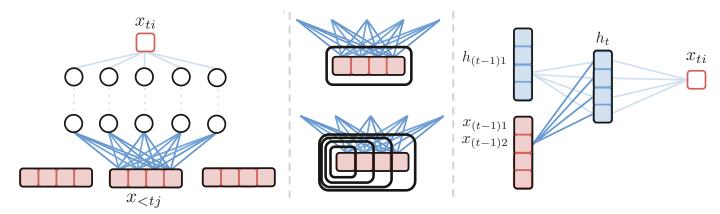


Fig. 1. (left) Schematic for modeling Granger causality using cMLPs. If the outgoing weights for series j, shown in dark blue, are penalized to zero, then series j does not Granger-cause series i. (center) The group lasso penalty jointly penalizes the full set of outgoing weights while the hierarchical version penalizes the nested set of outgoing weights, penalizing higher lags more. (right) Schematic for modeling Granger causality using a cLSTM. If the dark blue outgoing weights to the hidden units from an input  $x_{(t-1)j}$  are zero, then series j does not Granger-cause series i.

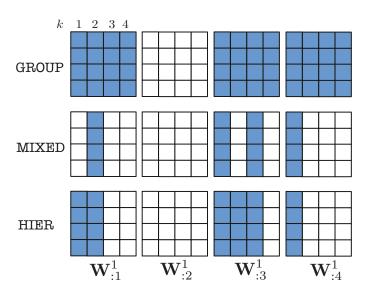


Fig. 2. Example of group sparsity patterns of the first layer weights of a cMLP with four first layer hidden units and four input series with maximum lag k = 4. Differing sparsity patterns are shown for the three different structured penalties of group lasso (GROUP) from Equation (9), group sparse group lasso (MIXED) from Equation (10) and hierarchical lasso (HIER) from Equation (11).

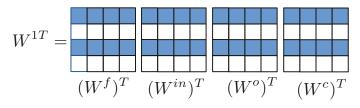


Fig. 3. Example of the group sparsity patterns in a sparse cLSTM model with a four dimensional hidden state and four input series. Due to the group lasso penalty on the columns of W, the  $W^f$ ,  $W^{in}$ ,  $W^o$ , and  $W^c$ matrices will share the same column sparsity pattern.

causality by applying a group penalty to the columns of the  $W^{1k}$  matrices for each  $q_i$ ,

$$\min_{\mathbf{W}} \sum_{t=K}^{T} \left( x_{it} - g_i \left( x_{(t-1):(t-K)} \right) \right)^2 + \lambda \sum_{j=1}^{p} \Omega \left( W_{:j}^1 \right). \quad (8)$$

where  $\Omega$  is a penalty that shrinks the entire set of first layer weights for input series j, i.e.,  $W_{:j}^1 = \left(W_{:j}^{11}, \dots, W_{:j}^{1K}\right)$ , to zero. We consider three different penalties that, together, show how we recast structured regression penalties to the neural network case.

We first consider a group lasso penalty over the entire set of outgoing weights across all lags for time series j,  $W_{ij}^1$ 

$$\Omega\left(W_{\cdot i}^{1}\right) = \left\|W_{\cdot i}^{1}\right\|_{F},\tag{9}$$

where  $\|\cdot\|_F$  is the Froebenius matrix norm. This penalty shrinks all weights associated with lags for input series j equally. For large enough  $\lambda$ , the solutions to Equation (8) with the group penalty in Equation (9) will lead to many zero columns in each  $W^{1k}$  matrix, implying only a small number of estimated Granger causal connections. This group penalty is the neural network analogue of the group lasso penalty across lags in Equation 2 for the VAR case.

To detect the lags where Granger causal effects exists, we propose a new penalty called a group sparse group lasso penalty. This penalty assumes that only a few lags of a series j are predictive of series i, and provides both sparsity across groups (a sparse set of Granger causal time series) and sparsity within groups (a subset of relevant lags)

$$\Omega\left(W_{:j}^{1}\right) = \alpha \left\|W_{:j}^{1}\right\|_{F} + (1 - \alpha) \sum_{k=1}^{K} \left\|W_{:j}^{1k}\right\|_{2}$$
 (10)

where  $\alpha \in (0,1)$  controls the tradeoff in sparsity across and within groups. This penalty is a related to, and is a generalization of, the sparse group lasso [49].

Finally, we may simultaneously select for both Granger causality and the lag order of the interaction by replacing the group lasso penalty in Equation (8) with a hierarchical group lasso penalty [15] in the MLP optimization problem,

$$\Omega\left(W_{:j}^{1}\right) = \sum_{k=1}^{K} \left\| \left(W_{:j}^{1k}, \dots, W_{:j}^{1K}\right) \right\|_{F}.$$
 (11)

The hierarchical penalty leads to solutions such that for each  $\min_{\mathbf{W}} \sum_{t=K}^{T} \left( x_{it} - g_i \left( x_{(t-1):(t-K)} \right) \right)^2 + \lambda \sum_{i=1}^{p} \Omega \left( W_{:j}^1 \right). \quad \text{(8)} \quad \text{if the rectation periods it is solutions as an extension of the periods of the period of the periods of the period of the periods of the period of$  the lag of each interaction. The hierarchical penalty also sets many columns of  $W^{1k}$  to be zero across all k, effectively selecting for Granger causality. In practice, the hierarchical penalty allows us to fix K to a large value, ensuring that no Granger causal connections at higher lags are missed. Example sparsity patterns selected by the three penalties are shown in Figure 2.

While the primary motivation of our penalties is for efficient Granger causality selection, the lag selection penalties in Equations (10) and (11) are also of independent interest to nonlinear forecasting with neural networks. In this case, over-specifying the lag of a NAR model leads to poor generalization and overfitting [25]. One proposed technique in the literature is to first select the appropriate lags using forward orthogonal least squares [25]; our approach instead combines model fitting and lag selection into one procedure.

#### 3.3 Sparse Input RNNs

Recurrent neural networks (RNNs) are particularly well suited for modeling time series, as they compress the past of a time series into a hidden state, aiming to capture complicated nonlinear dependencies at longer time lags than traditional time series models. As with MLPs, time series forecasting with RNNs typically proceeds by jointly modeling the entire evolution of the multivariate series using a single recurrent network.

As in the MLP case, it is difficult to disentangle how each series affects the evolution of another series when using an RNN. This problem is even more severe in complicated recurrent networks like LSTMs. To model Granger causality with RNNs, we follow the same strategy as with MLPs and model each  $g_i$  function using a separate RNN, which we refer to as a component-wise RNN (cRNN). For simplicity, we assume a single-layer RNN, but our formulation may be easily generalized to accommodate more layers.

Consider an RNN for predicting a single component. Let  $\mathbf{h}_t \in \mathbb{R}^H$  represent the H-dimensional hidden state at time t, representing the historical context of the time series for predicting a component  $x_{ti}$ . The hidden state at time t+1 is updated recursively

$$\mathbf{h}_t = f_i(\mathbf{x}_t, \mathbf{h}_{t-1}),\tag{12}$$

where  $f_i$  is some nonlinear function that depends on the particular recurrent architecture.

Due to their effectiveness at modeling complex time dependencies, we choose to model the recurrent function f using an LSTM [26]. The LSTM model introduces a second hidden state variable  $\mathbf{c}_t$ , referred to as the cell state, giving the full set of hidden parameter as  $(\mathbf{c}_t, \mathbf{h}_t)$ . The LSTM model updates its hidden states recursively as

$$\mathbf{f}_{t} = \sigma \left( W^{f} \mathbf{x}_{t} + U^{f} \mathbf{h}_{(t-1)} \right)$$

$$\mathbf{i}_{t} = \sigma \left( W^{in} \mathbf{x}_{t} + U^{in} \mathbf{h}_{(t-1)} \right)$$

$$\mathbf{o}_{t} = \sigma \left( W^{o} \mathbf{x}_{t} + U^{o} \mathbf{h}_{(t-1)} \right)$$

$$\mathbf{c}_{t} = \mathbf{f}_{t} \odot \mathbf{c}_{t-1} + \mathbf{i}_{t} \odot \sigma \left( W^{c} \mathbf{x}_{t} + U^{c} \mathbf{h}_{t-1} \right)$$

$$\mathbf{h}_{t} = \mathbf{o}_{t} \odot \sigma (\mathbf{c}_{t})$$

$$(13)$$

where  $\odot$  denotes element-wise multiplication and  $\mathbf{i}_t$ ,  $\mathbf{f}_t$ , and  $\mathbf{o}_t$  represent input, forget and output gates, respectively, that control how each component of the state cell,  $\mathbf{c}_t$ , is

updated and then transferred to the hidden state used for prediction,  $\mathbf{h}_t$ . In particular, the forget gate,  $\mathbf{f}_t$ , controls the amount that the past cell state influences the future cell state, and the input gate,  $\mathbf{i}_t$ , controls the amount that the current observation influences the new cell state. The additive form of the cell state update in the LSTM allows it to encode longrange dependencies, since cell states from far in the past may still influence the cell state at time t if the forget gates remain close to one. In the context of Granger causality, this flexible architecture can represent long-range, nonlinear dependencies between time series. As in the cMLP, the output for series i at time t is given by a linear decoding of the hidden state

$$x_{ti} = g_i(x_{< t}) + e_{ti} = W^2 \mathbf{h}_t + e_{ti},$$
 (14)

where  $W^2$  are the output weights. We let  $\mathbf{W} = (W^1, W^2, U^1)$  be the full set of parameters where  $W^1 = ((W^f)^\top, (W^{in})^\top, (W^o)^\top, (W^c)^\top)^\top$  and  $U^1 = ((U^f)^\top, (U^{in})^\top, (U^o)^\top, (U^c)^\top)^\top$  represent the full set of first layer weights. As in the MLP case, other decoding schemes could be used in the case of categorical or count data.

#### 3.3.1 Granger Causality Selection in LSTMs

In Equation (14) the set of input matrices  $W^1$  controls how the past time series affect the forget gates, input gates, output gates, and cell updates, and, consequently, the update of the hidden representation. Like in the MLP case, for this component-wise LSTM model (cLSTM) a sufficient condition for Granger non-causality of an input series j on an output i is that all elements of the jth column of  $W^1$  are zero,  $W^1_{:j}=0$ . Thus, we may select series that Granger-cause series i using a group lasso penalty across columns of  $W^1$  by

$$\min_{\mathbf{W}} \sum_{t=2}^{T} \left( x_{it} - g_i(x_{< t}) \right)^2 + \lambda \sum_{j=1}^{p} ||W_{:j}^1||_2.$$
 (15)

For a large enough  $\lambda$ , many columns of  $W^1$  will be zero, leading to a sparse set of Granger causal connections. An example sparsity pattern in the LSTM parameters is shown in Figure 3.

#### 4 OPTIMIZING THE PENALIZED OBJECTIVES

#### 4.1 Optimizing the Penalized cMLP Objective

We optimize the nonconvex objectives of Equation (8) using proximal gradient descent [50]. Proximal optimization is important in our context because it leads to exact zeros in the columns of the input matrices, a critical requirement for interpretating Granger non-causality in our framework. Additionally, a line search can be incorporated into the optimization algorithm to ensure convergence to a local minimum [51]. The algorithm updates the network weights  $\mathbf{W}$  iteratively starting with  $\mathbf{W}^{(0)}$  by

$$\mathbf{W}^{(m+1)} = \operatorname{prox}_{\gamma^{(m)}\lambda\Omega} \left( \mathbf{W}^{(m)} - \gamma^{(m)} \nabla \mathcal{L}(\mathbf{W}^{(m)}) \right), \quad (16)$$

where  $\mathcal{L} = \sum_{t=K}^{T} \left( x_{ti} - g_i(x_{< t}) \right)^2$  is the the neural network prediction loss and  $\operatorname{prox}_{\lambda\Omega}$  is the proximal operator with respect to the sparsity inducing penalty function  $\Omega$ . The

entries in  $\mathbf{W}^{(0)}$  are initialized randomly from a standard normal distribution. The scalar  $\gamma^{(m)}$  is the step size, which is either set to a fixed value or determined by line search [51]. While the objectives in Equation (8) are nonconvex, we find that no random restarts are required to accurately detect Granger causality connections.

Since the sparsity promoting group penalties are only applied to the input weights, the proximal step for weights at the higher levels is simply the identity function. The proximal step for the group lasso penalty on the input weights is given by a group soft-thresholding operation on the input weights [50],

$$\operatorname{prox}_{\gamma^{(m)}\lambda\Omega}(W_{:k}^{i}) = \operatorname{soft}(W_{:k}^{1}, \gamma^{(m)}\lambda) \tag{17}$$

$$= \left(1 - \frac{\lambda \gamma^{(m)}}{\|W_{:j}^1\|_F}\right)_+ W_{:k}^1, \qquad (18)$$

where  $(x)_+ = \max(0,x)$ . For the group sparse group lasso, the proximal step on the input weights is given by groupsoft thresholding on the lag specific weights, followed by group soft thresholding on the entire resulting input weights for each series, see Algorithm 2. The proximal step on the input weights for the hierarchical penalty is given by iteratively applying the group soft-thresholding operation on each nested group in the penalty, from the smallest group to the largest group [42], and is shown in Algorithm 3.

Algorithm 1 Proximal gradient descent with line search algorithm for solving Equation (8). Proximal steps given in Equation (17) for the group lasso penalty, in Algorithm 2 for the group sparse group lasso penalty, and in Algorithm 3 for the hierarchical penalty.

```
\begin{split} & \mathbf{Require:} \  \, \lambda > 0 \\ & m = 0 \text{, initialize } \mathbf{W}^{(0)} \\ & \mathbf{while} \text{ not converged } \mathbf{do} \\ & m = m+1 \\ & \text{ determine } \gamma \text{ by line search} \\ & \mathbf{for} \  \, j = 1 \text{ to } p \  \, \mathbf{do} \\ & W_{:j}^{1(m+1)} = \mathrm{prox}_{\gamma\lambda\Omega} \left( W_{:j}^{1(m)} - \gamma \nabla_{W_{:j}^1} \mathcal{L} \left( \mathbf{W}^{(m)} \right) \right) \\ & \mathbf{end for} \\ & \mathbf{for} \  \, l = 2 \text{ to } L \  \, \mathbf{do} \\ & W^{l(m+1)} = W^{l(m)} - \gamma \nabla_{W^l} \mathcal{L} \left( \mathbf{W}^{(m)} \right) \\ & \mathbf{end for} \\ & \mathbf{end for} \\ & \mathbf{end while} \\ & \mathbf{return} \  \, (\mathbf{W}^{(m)}) \end{split}
```

**Algorithm 2** One pass algorithm to compute the proximal map for the group sparse group lasso penalty, for relevant lag selection in the cMLP model.

**Algorithm 3** One pass algorithm to compute the proximal map for the hierarchical group lasso penalty, for automatic lag selection in the cMLP model.

$$\begin{array}{l} \textbf{Require:} \ \lambda > 0, \gamma > 0, \left(W^{11}_{:j}, \ldots, W^{1K}_{:j}\right) \\ \textbf{for} \ k = K \ \text{to} \ 1 \ \textbf{do} \\ \left(W^{1k}_{:j}, \ldots, W^{1K}_{:j}\right) = \operatorname{soft}\left(\left(W^{1k}_{:j}, \ldots, W^{1K}_{:j}\right), \gamma\lambda\right) \\ \textbf{end for} \\ \textbf{return} \ \left(W^{11}_{:j}, \ldots, W^{1K}_{:j}\right) \end{array}$$

Since all datasets we study are relatively small, the gradients are with respect to the full data objective (i.e., all time points); for larger datasets, one could instead use proximal stochastic gradient descent [52].

### 4.2 Optimizing the Penalized cLSTM Objective

Similar to the cMLP, we optimize Equation (15) using proximal gradient descent. When the data consists of many replicates of short time series, like in the DREAM3 data in Section 7, we perform a full backpropagation through time (BPTT) to compute the gradients. However, for longer series we truncate the BPTT by unlinking the hidden sequences. In practice, we do this by splitting the dataset up into equal sized batches, and treating each batch as an independent realization. Under this approach, the gradients used to optimize Equation (15) are only approximations of the gradients of the full cLSTM model. This is very common practice in the training of of RNNs [53], [54], [55]. The full optimization algorithm for training is shown in Algorithm 4.

**Algorithm 4** Proximal gradient descent with line search algorithm for solving Equation (15) for the cLSTM with group lasso penalty.

```
Require: \lambda > 0
m = 0, initialize \mathbf{W}^{(0)}
while not converged \mathbf{do}
m = m + 1
compute \nabla \mathcal{L}\left(\mathbf{W}^{(m)}\right) by BPTT (truncated for large T)
determine \gamma by line search.
for j = 1 to p \mathbf{do}
W_{:j}^{1(m+1)} = \operatorname{soft}\left(W_{:j}^{1(m)} - \gamma \nabla_{W_{:j}^{1}} \mathcal{L}\left(\mathbf{W}^{(m)}\right), \gamma \lambda\right)
end for
W^{2(m+1)} = W^{2(m)} - \gamma \nabla_{W^{2}} \mathcal{L}\left(\mathbf{W}^{(m)}\right)
U^{1(m+1)} = U^{1(m)} - \gamma \nabla_{U^{1}} \mathcal{L}\left(\mathbf{W}^{(m)}\right)
end while return (\mathbf{W}^{(m)})
```

# 5 COMPARING CMLP AND CLSTM MODELS FOR GRANGER CAUSALITY

Both the cMLP and cLSTM frameworks model each component function  $g_i$  using independent networks for each i. For the cMLP model, one needs to specify a maximum possible model lag K. However, our lag selection strategy (Equation 11) allows one to set K to a large value and the weights for higher lags are automatically removed from the model. On the other hand, the cLSTM model requires

TABLE 1

Comparison of AUROC for Granger causality selection among different approaches, as a function of the forcing constant F and the length of the time series T. Results are the mean across five initializations, with 95% confidence intervals.

Model		F = 10			F = 40	
$\overline{T}$	250	500	1000	250	500	1000
cMLP cLSTM IMV-LSTM LOO-LSTM	86.6 ± 0.2 81.3 ± 0.9 63.7 ± 4.3 47.9 ± 3.2	$96.6 \pm 0.2$ $93.4 \pm 0.7$ $76.0 \pm 4.5$ $49.4 \pm 1.8$	$98.4 \pm 0.1$ $96.0 \pm 0.1$ $85.5 \pm 3.4$ $50.1 \pm 1.0$	84.0 ± 0.5 75.1 ± 0.9 53.6 ± 5.2 50.1 ± 3.3	89.6 ± 0.2 87.8 ± 0.4 59.0 ± 4.5 49.1 ± 3.2	$95.5 \pm 0.3$ $94.4 \pm 0.5$ $69.0 \pm 4.8$ $51.1 \pm 3.7$

no maximum lag specification, and instead automatically learns the memory of each interaction. As a consequence, the cMLP and cLSTM differ in the amount of data used for training, as noted by a comparison of the t index in Equation (15) and Equation (11). For a length T series, the cMLP and cLSTM models use T-K and T-1 data points, respectively. While insignificant for large T, when the data consist of independent replicates of short series, as in the DREAM3 data in Section 7, the difference may be important. This ability to simultaneously model longer range dependencies while harnessing the full training set may explain the impressive performance of the cLSTM in the DREAM3 data in Section 7.

Finally, the zero outgoing weights in both the cMLP and cLSTM are a sufficient but not necessary condition to represent Granger non-causality. Indeed, series i could be Granger non-causal of series j through a complex configuration of weights that exactly cancel each other. However, because we wish to interpret the outgoing weights of the inputs as a measure of dependence, it is important that these weights reflect the true relationship between inputs and outputs. Our penalization schemes in both cMLP and cLSTM acts as a prior that biases the network to represent Granger non-causal relationships with zeros in the outgoing weights of the inputs, rather than through other configurations. Our simulation results in Section 6 validate this intuition.

#### 6 SIMULATION EXPERIMENTS

#### 6.1 cMLP and cLSTM Simulation Comparison

To compare and analyze the performance of our two approaches, the cMLP and cLSTM, we apply both methods to detecting Granger causality networks in simulated linear VAR data and simulated Lorenz-96 data [32], a nonlinear model of climate dynamics. Overall, the results show that our methods can accurately reconstruct the underlying Granger causality graph in both linear and nonlinear settings. We first describe the results from the Lorenz experiment and present the VAR results subsequently.

#### 6.1.1 Lorenz-96 Model

The continuous dynamics in a p-dimensional Lorenz model are given by

$$\frac{dx_{ti}}{dt} = (x_{t(i+1)} - x_{t(i-2)}) x_{t(i-1)} - x_{ti} + F, \quad (19)$$

where  $x_{t(-1)} = x_{t(p-1)}, x_{t0} = x_{tp}, x_{t(p+1)} = x_{t1}$  and F is a forcing constant that determines the level of nonlinearity and chaos in the series. Example series for two settings of F

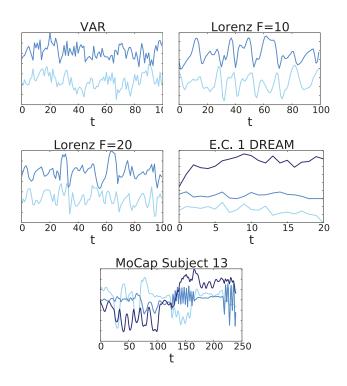


Fig. 4. Example multivariate linear (VAR) and nonlinear (Lorenz, DREAM, and MoCap) series that we analyze using both cMLP and cLSTM models. Note as the forcing constant, F, in the Lorenz model increases, the data become more chaotic.

are displayed in Figure 4. We numerically simulate a p=20 Lorenz-96 model with a sampling rate of  $\Delta_t=0.05$ , which results in a multivariate, nonlinear time series with sparse Granger causal connections.

Using this simulation setup, we test our models' ability to recover the underlying causal structure. Average values of area under the ROC curve (AUROC) for recovery of the causal structure across five initialization seeds are shown in Table 1, and we obtain results under three different data set lengths,  $T \in (250, 500, 1000)$ , and two forcing constants,  $F \in (10, 40)$ .

We compare results for the cMLP and cLSTM with two baseline methods that also rely on neural networks, the IMV-LSTM and leave-one-out LSTM (LOO-LSTM) approaches. The IMV-LSTM [56], [57] uses attention weights to provide greater interpretability than standard LSTMs, and it detects Granger causal relationships by aggregating its attention weights. LOO-LSTM detects Granger causal relationships through the increase in loss that results from

TABLE 2

Comparison of AUROC for Granger causality selection among different approaches, as a function of the VAR lag order and the length of the time series T. Results are the mean across five initializations, with 95% confidence intervals.

Model		VAR(1)			VAR(2)	
$\overline{T}$	250	500	1000	250	500	1000
cMLP cLSTM IMV-LSTM LOO-LSTM	91.6 ± 0.4 88.5 ± 0.9 53.7 ± 7.9 50.1 ± 2.7	$94.9 \pm 0.2$ $93.4 \pm 1.9$ $63.2 \pm 8.0$ $50.2 \pm 2.6$	$98.4 \pm 0.1$ $97.6 \pm 0.4$ $60.4 \pm 8.3$ $50.5 \pm 1.9$	$84.4 \pm 0.2$ $83.5 \pm 0.3$ $53.5 \pm 3.9$ $50.1 \pm 1.4$	$88.3 \pm 0.4$ $92.5 \pm 0.9$ $54.3 \pm 3.6$ $50.4 \pm 1.4$	$95.1 \pm 0.2$ $97.8 \pm 0.1$ $55.0 \pm 3.4$ $50.0 \pm 1.0$

withholding each input time series (see Appendix).

We use H=100 hidden units for all four methods, as experiments show that performance does not improve with a different number of units. While more layers may prove beneficial, for all experiments we fix the number of hidden layers, L, to one and leave the effects of additional hidden layers to future work. For the cMLP, we use the hierarchical penalty with model lag of K=5; see Section 6.2 for a performance comparison of several possible penalties across model input lags.

For our methods, we compute AUROC values by sweeping  $\lambda$  across a range of values; discarded edges (inferred Granger non-causality) for a particular  $\lambda$  setting are those whose associated  $L_2$  norm of the input weights of the neural network is equal to zero. Note that our proximal gradient algorithm sets many of these groups to be exactly zero. We compute AUROC values for the IMV-LSTM and LOO-LSTM by sweeping a range of thresholds for either the attention values or the increase in loss due to withholding time series.

As expected, the results indicate that the cMLP and cLSTM performance improves as the data set size T increases. The cMLP outperforms the cLSTM both in the less chaotic regime of F=10 and the more chaotic regime of F=40, but the gap in their performance narrows as more data is used. Both methods outperform the IMV-LSTM and LOO-LSTM by a wide margin. Our models' 95% confidence intervals are also relatively narrow, at less than 1% AUROC for the cLSTM and cMLP, compared with 3-5% for the IMV-LSTM.

To understand the role of the number of hidden units in our methods, we perform an ablation study to test different values of H; the results show that both the cMLP and cLSTM are robust to smaller H values, but that their performance benefits from H=100 hidden units (see Appendix). Additionally, we investigate the importance of the optimization algorithm; we found that Adam [58], proximal gradient descent [50] and proximal gradient descent with a line search [51] lead to similar results (see Appendix). However, because Adam requires a thresholding parameter and the line search is computationally costly, we use standard proximal gradient descent in the remainder of our experiments.

#### 6.1.2 VAR Model

To analyze our methods' performance when the true underlying dynamics are linear, we simulate data from p=20 VAR(1) and VAR(2) models with randomly generated sparse transition matrices. To generate sparse dependencies for each time series i, we create self dependencies and randomly select three more dependencies among the other

TABLE 3

AUROC comparisons between different cMLP Granger causality selection penalties on simulated Lorenz-96 data as a function of the input model lag, K. Results are the mean across five initializations, with 95% confidence intervals.

Lag K	5	10	20
GROUP	$88.1 \pm 0.8$	$82.5 \pm 0.3$	$80.5 \pm 0.5$
MIXED	$90.1 \pm 0.5$	$85.4 \pm 0.3$	$83.3 \pm 1.1$
HIER	$95.5 \pm 0.2$	$95.4 \pm 0.5$	$95.2 \pm 0.3$

p-1 time series. Where series i depends on series j, we set  $A_{ij}^k=0.1$  for k=1 or k=1,2, and all other entries of A are set to zero. Examining both VAR models allows us to see how well our methods detect Granger causality at longer time lags, even though no time lag is explicitly specified in our models. Our results are the average over five random initializations for a single dependency graph.

The AUROC results are displayed in Table 2 for the cLSTM, cMLP, IMV-LSTM, and LOO-LSTM approaches for three dataset lengths,  $T \in (250, 500, 1000)$ . The performance of the cLSTM and cMLP improves at larger T, and, as in the Lorenz-96 case, both models outperform the IMV-LSTM and LOO-LSTM by a wide margin. The cMLP remains more robust than the cLSTM with smaller amounts of data, but the cLSTM outperforms the cMLP on several occasions with T=500 or T=1000.

The IMV-LSTM consistently underperforms our methods with these datasets, likely because it is not explicitly designed for Granger causality discovery. Our finding that the IMV-LSTM performs poorly at this task is consistent with recent work suggesting that attention mechanisms are not indicative of feature importance [59]. The LOO-LSTM approach consistently achieves poor performance, likely due to two factors: (i) unregularized LSTMs are prone to overfitting in the low-data regime, even when Granger causal time series are held out, and (ii) withholding a single time series will not impact the loss if the remaining time series have dependencies that retain its signal.

#### 6.2 Quantitative Analysis of the Hierarchical Penalty

We next quantitatively compare the three possible structured penalties for Granger causality selection in the cMLP model. In Section 3.2 we introduced the full group lasso (GROUP) penalty over all lags (Equation 8), the group sparse group lasso (MIXED) (Equation 10) and the hierarchical (HIER) lag selection penalty (Equation 11). We compare these approaches across various choices of the cMLP model's maximum lag,  $K \in (5, 10, 20)$ . We use H = 10

hidden units for data simulated from the nonlinear Lorenz model with F=20, p=20, and T=750. As in Section 6.1, we compute the mean AUROC over five random initializations and display the results in Table 3. Importantly, the hierarchical penalty outperforms both group and mixed penalties across all model input lags K. Furthermore, performance significantly declines as K increases in both group and mixed settings while the performance of the hierarchical penalty stays roughly constant as K increases. This result suggests that performance of the hierarchical penalty for nonlinear Granger causality selection is robust to the input lag, implying that precise lag specification is unnecessary. In practice, this allows one to set the model lag to a large value without worrying that nonlinear Granger causality detection will be compromised.

#### 6.3 Qualitative Analysis of the Hierarchical Penalty

To qualitatively validate the performance of the hierarchical group lasso penalty for automatic lag selection, we apply our penalized cMLP framework to data generated from a sparse VAR model with longer interactions. Specifically, we generate data from a p=10, VAR(3) model as in Section 2. To generate sparse dependencies for each time series i, we create self dependencies and randomly select two more dependencies among the other p-1 time series. When series i depends on series j, we set  $A^k_{ij}=0.1$  for k=1,2,3. All other entries of A are set to zero. This implies that the Granger causal connections that do exist are of true lag 3. We run the cMLP with the hierarchical group lasso penalty and a maximal lag order of K=5; for comparison, we also train a VAR model with a hierarchical penalty and maximal lag order K=5.

We visually display the selection results for one cMLP (i.e., one output series) and the VAR baseline across a variety of  $\lambda$  settings in Figure 5. For the lower  $\lambda=4.09e-4$  setting, the cMLP both (i) overestimates the lag order for a few input series and (ii) allows some false positive Granger causal connections. For the higher  $\lambda=7.94e-4$ , lag selection performs almost perfectly, in addition to correct estimation of the Granger causality graph. Higher  $\lambda$  values lead to larger penalization on longer lags, resulting in weaker long-lag connections. The VAR model, which is ideal for VAR data, does not perform noticeably better. While we show results for multiple  $\lambda$  values for visualization, in practice one may use cross validation to select the appropriate  $\lambda$ .

#### 7 DREAM CHALLENGE

We next apply our methods to estimate Granger causality networks from a realistically simulated time course gene expression data set. The data are from the DREAM3 challenge [35] and provide a difficult, nonlinear data set for rigorously comparing Granger causality detection methods [33], [34]. The data is simulated using continuous gene expression and regulation dynamics, with multiple hidden factors that are not observed. The challenge contains five different simulated data sets, each with different ground truth Granger causality graphs: two E. Coli (E.C.) data sets and three Yeast (Y.) data sets. Each data set contains p=100 different time series, each with 46 replicates sampled at 21

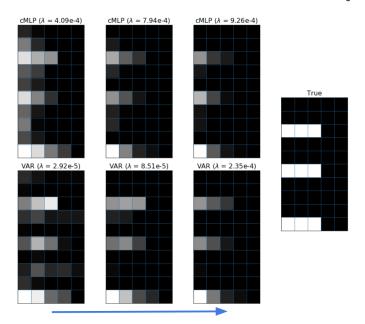


Fig. 5. Qualitative results of the cMLP automatic lag selection using a hierarchical group lasso penalty and maximal lag of K=5. The true data are from a VAR(3) model. The images display results for a single cMLP (one output series) and a VAR model using various penalty strengths  $\lambda$ . The rows of each image correspond to different input series while the columns correspond to the lag, with k=1 at the left and k=5 at the right. The magnitude of each entry is the  $L_2$  norm of the associated input weights of the neural network after training. The true lag interactions are shown in the rightmost image. Brighter color represents larger magnitude.

time points for a total of 966 time points. This represents a very limited data scenario relative to the dimensionality of the networks and complexity of the underlying dynamics of interaction. Three time series components from a single replicate of the E. Coli 1 data set are shown in Figure 4.

We apply both the cMLP and cLSTM to all five data sets. Due to the short length of the series replicates, we choose the maximum lag in the cMLP to be K=2 and use H=10 and H=5 hidden units for the cMLP and cLSTM, respectively. For our performance metric, we consider the DREAM3 challenge metrics of area under the ROC curve (AUROC) and area under the precision recall curve (AUPR). Both curves are computed by sweeping  $\lambda$  over a range of values, as described in Section 6.

In Figure 6, we compare the AUROC and AUPR of our cMLP and cLSTM to previously published AUROC and AUPR results on the DREAM3 data [33]. These comparisons include both linear and nonlinear approaches: (i) a linear VAR model with a lasso penalty (LASSO) [7], (ii) a dynamic Bayesian network using first-order conditional dependencies (G1DBN) [34], and (iii) a state-of-the-art multioutput kernel regression method (OKVAR) [33]. The latter is the most mature of a sequence of nonlinear kernel Granger causality detection methods [60], [61]. In terms of AUROC, our cLSTM outperforms all methods across all five datasets. Furthermore, the cMLP method outperforms previous methods on two datasets, Y.1 and Y.3, ties G1DBN on Y.2, and slightly under performs OKVAR in E.C.1 and E.C.2. In terms of AUPR, both cLSTM and cMLP methods do much better than all previous approaches, with the cLSTM

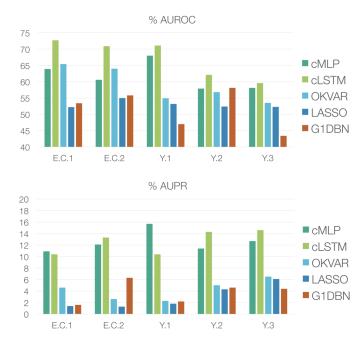


Fig. 6. (Top) AUROC and (bottom) AUPR (given in %) results for our proposed regularized cMLP and cLSTM models and the set of methods—OKVAR, LASSO, and G1DBN presented in [33]. These results are for the DREAM3 size-100 networks using the original DREAM3 data sets.

outperforming the cMLP in three datasets. The raw ROC curves for cMLP and cLSTM are displayed in Figure 7.

These results clearly demonstrate the importance of taking a nonlinear approach to Granger causality detection in a (simulated) real-world scenario. Among the nonlinear approaches, the neural network methods are extremely powerful. Furthermore, the cLSTM's ability to efficiently capture long memory (without relying on long-lag specifications) appears to be particularly useful. This result validates many findings in the literature where LSTMs outperform autoregressive MLPs. An interesting facet of these results, however, is that the impressive performance gains are achieved in a limited data scenario and on a task where the goal is recovery of interpretable structure. This is in contrast to the standard story of prediction on large datasets. To achieve these results, the regularization and induced sparsity of our penalties is critical.

## 8 DEPENDENCIES IN HUMAN MOTION CAPTURE DATA

We next apply our methodology to detect complex, nonlinear dependencies in human motion capture (MoCap) recordings. In contrast to the DREAM3 challenge results, this analysis allows us to more easily visualize and interpret the learned network. Human motion has been previously modeled using both linear dynamical systems [62], switching linear dynamical systems [37], [63] and also nonlinear dynamical models using Gaussian processes [64]. While the focus of previous work has been on motion classification [62] and segmentation [37], our analysis delves into the potentially long-range, nonlinear dependencies between different regions of the body during natural motion behavior. We consider a data set from the CMU MoCap database

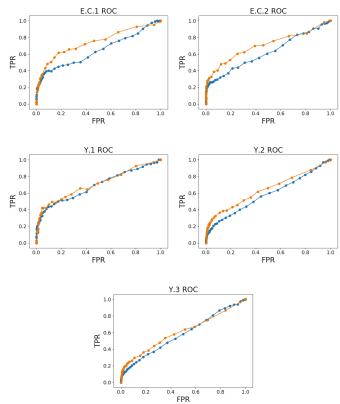


Fig. 7. ROC curves for the cMLP ( —— ) and cLSTM ( —— ) models on the five DREAM datasets.

[36] previously studied in [37]. The data set consists of p=54 joint angle and body position recordings across two different subjects for a total of T=2024 time points. In total, there are recordings from 24 unique regions because some regions, like the thorax, contain multiple angles of motion corresponding to the degrees of freedom of that part of the body.

We apply the cLSTM model with H=8 hidden units to this data set. For computational speed ups, we break the original series into length 20 segments and fit the penalized cLSTM model from Equation (15) over a range of  $\lambda$  values. To develop a weighted graph for visualization, we let the edge weight  $w_{ij}$  between components be the norm of the outgoing cLSTM weights from input series j to output component series i, standardized by the maximum such edge weight associated with the cLSTM for series i. Edges associated with more than one degrees of freedom (angle directions for the same body part) are averaged together. Finally, to aid visualization, we further threshold edge weights of magnitude 0.01 and below.

The resulting estimated graphs are displayed in Figure 9 for multiple values of the regularization parameter,  $\lambda$ . While we present the results for multiple  $\lambda$ , one may use cross validation to select  $\lambda$  if one graph is required. To interpret the presented skeleton plots, it is useful to understand the full set of motion behaviors exhibited in this data set. These behaviors are depicted in Figure 8, and include instances of *jumping jacks*, *side twists*, *arm circles*, *knee raises*, *squats*, *punching*, various forms of *toe touches*, and *running in place*. Due to the extremely limited data for any individual be-

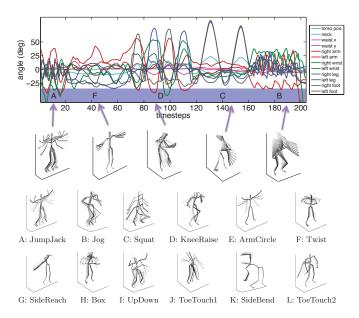


Fig. 8. (Top) Example time series from the MoCap data set paired with their particular motion behaviors. (Bottom) Skeleton visualizations of 12 possible exercise behavior types observed across all sequences analyzed in the main text.

havior, we chose to learn interactions from data aggregated over the entire collection of behaviors. In Figure 9, we see many intuitive learned interactions. For example, even in the more sparse graph (largest  $\lambda$ ) we learn a directed edge from right knee to left knee and a separate edge from left knee to right. This makes sense as most human motion, including the motions in this dataset involving lower body movement, entail the right knee leading the left and then vice versa. We also see directed interactions leading down each arm, and between the hands and toes for toe touches.

#### 9 CONCLUSION

We have presented a framework for nonlinear Granger causality selection using regularized neural network models of time series. To disentangle the effects of the past of an input series on the future of an output, we model each output series using a separate neural network. We then apply both the component multilayer perceptron (cMLP) and component long-short term memory (cLSTM) architectures, with associated sparsity promoting penalties on incoming weights to the network, and select for Granger causality. Overall, our results show that these methods outperform existing Granger causality approaches on the challenging DREAM3 data set and discover interpretable and insightful structure on a human MoCap data set.

Our work opens the door to multiple exciting avenues for future work. While we are the first to use a hierarchical lasso penalty in a neural network, it would be interesting to also explore other types of structured penalties, such as tree structured penalties [31].

Furthermore, although we have presented two relatively simple approaches, based off MLPs and LSTMs, our general framework of penalized input weights easily accommodates more powerful architectures. Exploring the effects of multiple hidden layers, powerful recurrent and convolutional

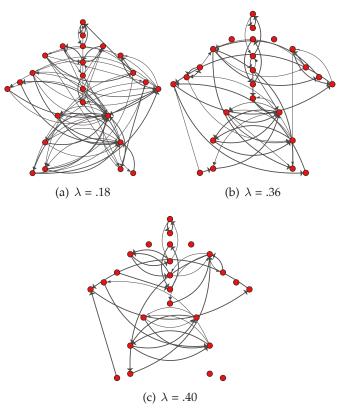


Fig. 9. Nonlinear Granger causality graphs inferred from the human Mo-Cap data set using the regularized cLSTM model. Results are displayed for a range of  $\lambda$  values. Each node corresponds to one location on the body.

architectures, like clockwork RNNs [65], and dilated causal convolutions [66], open up a wide range of research directions and the potential to detect long-range and complex dependencies. Further theoretical work on the identifiability of Granger non-causality in these more complex network models becomes even more important.

Finally, while we consider sparse input models, a different sparse output architecture would use a network, like an RNN, to learn hidden representations of each individual input series, and then model each output component as a sparse nonlinear combination across the hidden states of all time series, allowing a shared hidden representation across component tasks. A schematic of the proposed architecture that combines ideas from our cMLP and cLSTM models is shown in Figure 10.

## APPENDIX A MODEL ABLATIONS

We ran two ablation studies to understand factors that influence our methods' performance. First, we tested the cMLP and cLSTM with different numbers of hidden units on the Lorenz-96 data. Table 4 shows the AUROC results from a single run for two datasets with forcing constants  $F \in (10,40)$  and time series length T=1000, using different numbers of hidden units,  $H \in (5,10,25,50,100)$ . The results reveal that both models are robust to a small number of hidden units, but that their performance improves with

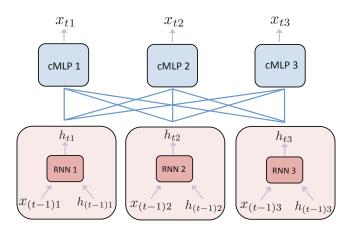


Fig. 10. Proposed architecture for detecting nonlinear Granger causality that combines aspects of both the cLSTM and cMLP models. A separate hidden representation,  $h_{tj}$ , is learned for each series j using an RNN. At each time point, the hidden states are each fed into a sparse cMLP to predict the individual output for each series  $x_{ti}$ . Joint learning of the whole network with a group penalty on the input weights of the individual cMLPs would allow the network to share information about hidden features in each  $h_{tj}$  while also allowing interpretable structure learning between the hidden states of each series and each output.

larger values of H. These findings suggest that overparameterization can help with the nonconvex optimization objective, leading to solutions that achieve high predictive accuracy while minimizing the penalty from the sparsity-inducing regularizer.

Next, we tested three approaches for optimizing our penalized objectives (Equations 8 and 15). We compared standard gradient descent with Adam [58] to proximal gradient descent (ISTA) [50] and proximal gradient descent with a line search (GIST) [51] on the Lorenz-96 data with T=1000 time points. Table 5 displays AUROC results across five initializations for two forcing constants  $F\in(10,40)$ , using the cMLP with H=10 hidden units. The results show that the three methods lead to similar results for both F=10 and F=40, although we did not compare the optimizers in other scenarios, e.g., with lower T values or with the cLSTM.

Among these optimization approaches, Adam is fastest due to its adaptive learning rate, but it requires a parameter for thresholding the resulting weights (while the proximal methods lead to exact zeros). In contrast, GIST guarantees convergence to a local minimum and is less sensitive to the learning rate parameter, but it is also considerably slower than Adam and ISTA. We therefore use standard proximal gradient descent, or ISTA, in the remainder of our experiments, because it leads to exact zeros while being more efficient than GIST. In practice, this means running Algorithm 1 or Algorithm 4 using a fixed learning rate  $\gamma$  rather than determining it by a line search.

## APPENDIX B BASELINE METHODS

The IMV-LSTM uses an attention mechanism to highlight the model's dependence on different parts of the input [57]. We train a separate IMV-LSTM model to predict each time series using all the time series as inputs, using the "IMV-Full" variant [57], and we use the attention weights from

TABLE 4 AUROC comparisons for the cMLP and cLSTM as a function of the number of hidden units  ${\cal H}$  for simulated Lorenz-96 data. Results are calculated using a single run.

Model	cMLP		cLSTM	
$\overline{F}$	10	40	10	40
H = 5 H = 10 H = 25 H = 50 H = 100	96.5 98.0 98.4 98.3 98.5	91.0 94.0 94.3 94.4 94.5	91.9 94.5 95.6 95.7 95.7	86.9 91.5 92.3 93.8 95.2

TABLE 5
AUROC comparisons between different optimization approaches for the cMLP with simulated Lorenz-96 data. Results are the mean across five initializations, with 95% confidence intervals.

F	10	40
GISTA	$98.0 \pm 0.2$	$93.8 \pm 0.3$
ISTA	$98.0 \pm 0.2$	$94.1 \pm 1.9$
Adam	$98.3 \pm 0.1$	$95.1 \pm 0.2$

the trained models to infer Granger causal relationships. Similar to the original work [56], we record the empirical mean of the attention values for each input time series for each model, and we construct a  $p \times p$  matrix of these values for the separate IMV-LSTMs. We then sweep over a range of threshold values to determine the most influential inputs for each IMV-LSTM, and we trace out an ROC curve from which we calculate AUROC values.

The LOO-LSTM baseline is based on the idea that withholding a highly predictive input should result in a decrease in predictive accuracy, a direction that has been explored for providing model-agnostic notions of feature importance [67], [68]. We begin by training separate LSTM models to predict each time series using all time series as inputs. We then train separate LSTM models to predict each time series i using all inputs except time series j, and we record the increase in loss when the jth time series is withheld. Using the results, we construct a  $p \times p$  matrix representing the differences in the loss, we sweep over a range of threshold values to determine the most influential inputs for each time series, and we trace out an ROC curve from which we calculate AUROC values.

#### **ACKNOWLEDGMENTS**

AT, IC, NF and EF acknowledge the support of ONR Grant N00014-15-1-2380, NSF CAREER Award IIS-1350133, and AFOSR Grant FA9550-16-1-0038. AT and AS acknowledge the support from NSF grants DMS-1161565 and DMS-1561814 and NIH grants R01GM114029 and R01GM133848.

### REFERENCES

- [1] O. Sporns, Networks of the Brain. MIT Press, 2010.
- [2] R. Vicente, M. Wibral, M. Lindner, and G. Pipa, "Transfer entropy—a model-free measure of effective connectivity for the neurosciences," *Journal of Computational Neuroscience*, vol. 30, no. 1, pp. 45–67, 2011.

- [3] P. A. Stokes and P. L. Purdon, "A study of problems encountered in Granger causality analysis from a neuroscience perspective," *Proceedings of the National Academy of Sciences*, vol. 114, no. 34, pp. E7063–E7072, 2017.
- [4] A. Sheikhattar, S. Miran, J. Liu, J. B. Fritz, S. A. Shamma, P. O. Kanold, and B. Babadi, "Extracting neuronal functional network dynamics via adaptive Granger causality analysis," *Proceedings of the National Academy of Sciences*, vol. 115, no. 17, pp. E3869–E3878, 2018.
- [5] W. F. Sharpe, G. J. Alexander, and J. W. Bailey, *Investments*. Prentice Hall, 1968.
- [6] A. Fujita, P. Severino, J. R. Sato, and S. Miyano, "Granger causality in systems biology: modeling gene networks in time series microarray data using vector autoregressive models," in *Brazilian Symposium on Bioinformatics*. Springer, 2010, pp. 13–24.
- [7] A. C. Lozano, N. Abe, Y. Liu, and S. Rosset, "Grouped graphical granger modeling methods for temporal causal modeling," in Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. ACM, 2009, pp. 577–586.
- [8] C. W. Granger, "Investigating causal relations by econometric models and cross-spectral methods," *Econometrica: Journal of the Econometric Society*, pp. 424–438, 1969.
- [9] H. Lütkepohl, New Introduction to Multiple Time Series Analysis. Springer Science & Business Media, 2005.
- [10] S. Basu, A. Shojaie, and G. Michailidis, "Network Granger causality with inherent grouping structure," The Journal of Machine Learning Research, vol. 16, no. 1, pp. 417–453, 2015.
- [11] K. Sameshima and L. A. Baccala, Methods in Brain Connectivity Inference Through Multivariate Time Series Analysis. CRC Press, 2016.
- [12] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.
- [13] M. Yuan and Y. Lin, "Model selection and estimation in regression with grouped variables," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 68, no. 1, pp. 49–67, 2006.
- [14] S. Basu, G. Michailidis *et al.*, "Regularized estimation in sparse high-dimensional time series models," *The Annals of Statistics*, vol. 43, no. 4, pp. 1535–1567, 2015.
- [15] W. B. Nicholson, J. Bien, and D. S. Matteson, "Hierarchical vector autoregression," arXiv preprint arXiv:1412.5250, 2014.
- [16] A. Shojaie and G. Michailidis, "Discovering graphical Granger causality using the truncating lasso penalty," *Bioinformatics*, vol. 26, no. 18, pp. i517–i523, 2010.
- [17] P. Bühlmann and S. Van De Geer, Statistics for High-Dimensional Data: Methods, Theory and Applications. Springer Science & Business Media, 2011.
- [18] T. Teräsvirta, D. Tjøstheim, C. W. J. Granger et al., Modelling Nonlinear Economic Time Series. Oxford University Press Oxford, 2010.
- [19] H. Tong, "Nonlinear time series analysis," International Encyclopedia of Statistical Science, pp. 955–958, 2011.
- [20] B. Lusch, P. D. Maia, and J. N. Kutz, "Inferring connectivity in networked dynamical systems: Challenges using Granger causality," *Physical Review E*, vol. 94, no. 3, p. 032220, 2016.
- [21] P.-O. Amblard and O. J. Michel, "On directed information theory and Granger causality graphs," *Journal of Computational Neuro*science, vol. 30, no. 1, pp. 7–16, 2011.
- [22] J. Runge, J. Heitzig, V. Petoukhov, and J. Kurths, "Escaping the curse of dimensionality in estimating multivariate transfer entropy," *Physical Review Letters*, vol. 108, no. 25, p. 258701, 2012.
- [23] M. Raissi, P. Perdikaris, and G. E. Karniadakis, "Multistep neural networks for data-driven discovery of nonlinear dynamical systems," arXiv preprint arXiv:1801.01236, 2018.
- [24] Ö. Kişi, "River flow modeling using artificial neural networks," Journal of Hydrologic Engineering, vol. 9, no. 1, pp. 60–63, 2004.
- [25] S. A. Billings, Nonlinear System Identification: NARMAX Methods in the Time, Frequency, and Spatio-Temporal Domains. John Wiley & Sons, 2013.
- [26] A. Graves, "Supervised sequence labelling," in Supervised Sequence Labelling with Recurrent Neural Networks. Springer, 2012, pp. 5–13.
- [27] R. Yu, S. Zheng, A. Anandkumar, and Y. Yue, "Long-term forecasting using tensor-train RNNs," arXiv preprint arXiv:1711.00073, 2017.
- [28] G. P. Zhang, "Time series forecasting using a hybrid arima and neural network model," *Neurocomputing*, vol. 50, pp. 159–175, 2003

- [29] Y. Li, R. Yu, C. Shahabi, and Y. Liu, "Graph convolutional recurrent neural network: Data-driven traffic forecasting," *arXiv preprint arXiv:1707.01926*, 2017.
- [30] J. Huang, T. Zhang, and D. Metaxas, "Learning with structured sparsity," *Journal of Machine Learning Research*, vol. 12, no. Nov, pp. 3371–3412, 2011.
- [31] S. Kim and E. P. Xing, "Tree-guided group lasso for multi-task regression with structured sparsity." in *International Conference on Machine Learning*, vol. 2. Citeseer, 2010, p. 1.
- [32] A. Karimi and M. R. Paul, "Extensive chaos in the lorenz-96 model," Chaos: An Interdisciplinary Journal of Nonlinear Science, vol. 20, no. 4, p. 043105, 2010.
- [33] N. Lim, F. D'Alché-Buc, C. Auliac, and G. Michailidis, "Operatorvalued kernel-based vector autoregressive models for network inference," *Machine Learning*, vol. 99, no. 3, pp. 489–513, 2015.
- [34] S. Lèbre, "Inferring dynamic genetic networks with low order independencies," Statistical Applications in Genetics and Molecular Biology, vol. 8, no. 1, pp. 1–38, 2009.
- [35] R. J. Prill, D. Marbach, J. Saez-Rodriguez, P. K. Sorger, L. G. Alexopoulos, X. Xue, N. D. Clarke, G. Altan-Bonnet, and G. Stolovitzky, "Towards a rigorous assessment of systems biology models: the dream3 challenges," *PloS One*, vol. 5, no. 2, p. e9202, 2010.
- [36] CMU, "Carnegie mellon university motion capture database," 2009, data retrieved from CMU, /http://mocap.cs.cmu.edu/.
- [37] E. B. Fox, M. C. Hughes, E. B. Sudderth, M. I. Jordan *et al.*, "Joint modeling of multiple time series via the beta process with application to motion capture segmentation," *The Annals of Applied Statistics*, vol. 8, no. 3, pp. 1281–1313, 2014.
- [38] J. M. Alvarez and M. Salzmann, "Learning the number of neurons in deep networks," in Advances in Neural Information Processing Systems, 2016, pp. 2270–2278.
- [39] C. Louizos, K. Ullrich, and M. Welling, "Bayesian compression for deep learning," in Advances in Neural Information Processing Systems, 2017, pp. 3288–3298.
- [40] J. Feng and N. Simon, "Sparse-input neural networks for highdimensional nonparametric regression and classification," arXiv preprint arXiv:1711.07592, 2017.
- [41] A. C. Lozano, N. Abe, Y. Liu, and S. Rosset, "Grouped graphical granger modeling for gene expression regulatory networks discovery," *Bioinformatics*, vol. 25, no. 12, pp. i110–i118, 2009.
- [42] R. Jenatton, J. Mairal, G. Obozinski, and F. Bach, "Proximal methods for hierarchical sparse coding," *Journal of Machine Learning Research*, vol. 12, no. Jul, pp. 2297–2334, 2011.
- [43] S. R. Chu, R. Shoureshi, and M. Tenorio, "Neural networks for system identification," *IEEE Control Systems Magazine*, vol. 10, no. 3, pp. 31–35, 1990.
- [44] S. Billings and S. Chen, "The determination of multivariable nonlinear models for dynamic systems using neural networks," 1996.
- [45] Y. Tao, L. Ma, W. Zhang, J. Liu, W. Liu, and Q. Du, "Hierarchical attention-based recurrent highway networks for time series prediction," arXiv preprint arXiv:1806.00685, 2018.
- [46] P. McCullagh and J. A. Nelder, Generalized Linear Models. CRC Press, 1989, vol. 37.
- [47] E. C. Hall, G. Raskutti, and R. Willett, "Inference of highdimensional autoregressive generalized linear models," arXiv preprint arXiv:1605.02693, 2016.
- [48] A. Tank, E. B. Fox, and A. Shojaie, "Granger causality networks for categorical time series," arXiv preprint arXiv:1706.02781, 2017.
- [49] N. Simon, J. Friedman, T. Hastie, and R. Tibshirani, "A sparse-group lasso," *Journal of Computational and Graphical Statistics*, vol. 22, no. 2, pp. 231–245, 2013. [Online]. Available: https://doi.org/10.1080/10618600.2012.681250
- [50] N. Parikh, S. Boyd *et al.*, "Proximal algorithms," Foundations and Trends in Optimization, vol. 1, no. 3, pp. 127–239, 2014.
- [51] P. Gong, C. Zhang, Z. Lu, J. Huang, and J. Ye, "A general iterative shrinkage and thresholding algorithm for non-convex regularized optimization problems," in *International Conference on Machine Learning*. PMLR, 2013, pp. 37–45.
- [52] L. Xiao and T. Zhang, "A proximal stochastic gradient method with progressive variance reduction," SIAM Journal on Optimization, vol. 24, no. 4, pp. 2057–2075, 2014.
- [53] R. J. Williams and D. Zipser, "Gradient-based learning algorithms for recurrent," *Backpropagation: Theory, architectures, and applica*tions, vol. 433, 1995.
- [54] I. Sutskever, Training Recurrent Neural Networks. University of Toronto Toronto, Canada, 2013.

- [55] P. J. Werbos et al., "Backpropagation through time: what it does and how to do it," Proceedings of the IEEE, vol. 78, no. 10, pp. 1550– 1560, 1990.
- [56] T. Guo, T. Lin, and Y. Lu, "An interpretable LSTM neural network for autoregressive exogenous model," arXiv preprint arXiv:1804.05251, 2018.
- [57] T. Guo, T. Lin, and N. Antulov-Fantulin, "Exploring interpretable LSTM neural networks over multi-variable data," in *International Conference on Machine Learning*. PMLR, 2019, pp. 2494–2504.
- [58] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," arXiv preprint arXiv:1412.6980, 2014.
- [59] S. Wiegreffe and Y. Pinter, "Attention is not not explanation," arXiv preprint arXiv:1908.04626, 2019.
- [60] V. Sindhwani, M. H. Quang, and A. C. Lozano, "Scalable matrix-valued kernel learning for high-dimensional nonlinear multivariate regression and Granger causality," arXiv preprint arXiv:1210.4792, 2012.
- [61] D. Marinazzo, M. Pellicoro, and S. Stramaglia, "Kernel-Granger causality and the analysis of dynamical networks," *Physical Review E*, vol. 77, no. 5, p. 056215, 2008.
- [62] E. Hsu, K. Pulli, and J. Popović, "Style translation for human motion," ACM Trans. Graph., vol. 24, no. 3, pp. 1082–1089, Jul. 2005. [Online]. Available: http://doi.acm.org/10.1145/1073204.1073315
- [63] V. Pavlovic, J. M. Rehg, and J. MacCormick, "Learning switching linear models of human motion," in *Advances in Neural Information Processing Systems*, 2001, pp. 981–987.
- [64] J. M. Wang, D. J. Fleet, and A. Hertzmann, "Gaussian process dynamical models for human motion," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 30, no. 2, pp. 283–298, 2007.
- [65] J. Koutnik, K. Greff, F. Gomez, and J. Schmidhuber, "A clockwork rnn," arXiv preprint arXiv:1402.3511, 2014.
- [66] A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, "Wavenet: A generative model for raw audio," arXiv preprint arXiv:1609.03499, 2016.
- [67] J. Lei, M. G'Sell, A. Rinaldo, R. J. Tibshirani, and L. Wasserman, "Distribution-free predictive inference for regression," *Journal of the American Statistical Association*, vol. 113, no. 523, pp. 1094–1111, 2018
- [68] G. Hooker and L. Mentch, "Please stop permuting features: An explanation and alternatives," arXiv preprint arXiv:1905.03151, 2019.