

# Granger-Causal Link Discovery in Large Temporal Networks through Conditional Models

Ananth Balashankar<sup>a</sup>, Srikanth Jagabathula<sup>b</sup>, and Lakshminarayanan Subramanian<sup>a,\*</sup>

<sup>a</sup>Department of Computer Science, New York University; <sup>b</sup>NYU Stern School of Business

This manuscript was compiled on January 11, 2022

Granger-causality derived from observational time series data is used in many real-world applications where timely interventions are infeasible. However, discovering Granger-causal links in large temporal networks with a large number of nodes and time-lags can lead to millions of time-lagged model parameters, which requires us to make sparsity and overlap assumptions. In this paper, we propose to learn time-lagged model parameters with the objective of improving recall of links, while learning to defer predictions when the overlap assumption is violated over observed time series. By learning such conditional time-lagged models, we demonstrate a 25% increase in the area under the precision-recall curve for discovering Granger-causal links combined with a 18-25% improvement in forecasting accuracy across three popular and diverse datasets from different disciplines (DREAM3 gene expression, MoCAP human motion recognition and New York Times news-based stock price prediction) with correspondingly large temporal networks, over several baseline models including Multivariate Autoregression, Neural Granger Causality, Graph Neural Networks and Graph Attention models. The observed improvement in Granger-causal link discovery is significant and can potentially further improve prediction accuracy and modeling efficiency in downstream real-world applications leveraging these popular datasets.

time series | granger causality | predictive models | overlap assumptions

## 1. Introduction

Granger causality (1) in time series data is important in many real world applications in economics (2), climate science (3) and biology (4). The knowledge of the Granger-causal structure allows us to build prediction models to make fine-grained time series forecasts conditioned on specific covariate values. For instance, the knowledge that a specific set of genes interfere with the expression of another gene allows us to build accurate gene regulatory networks, which assist in generating hypotheses for drug discovery. In practice, interventions are often infeasible because of the large dimensionality of data and making inference from real-time observations becomes inevitable as placing controls are either impractical, or even unethical. As a result, the inferences and forecasts must be done using only observational data. In this work, we assume that the underlying Granger-causal structure is specified to the extent that we know which covariates affect the outcome variable of interest. However, we lack information on how long the effect lasts and if it holds under all covariate distributions. For example, in the DREAM3 gene expression network task (5), where multiple genes can express and influence other genes, experiments are carried out where as part of a treatment, certain catalyzing agents are introduced over a time period and the corresponding gene expression time series are observed. Here, although we know that there are a specific set of 100 genes that potentially Granger-cause each other, we

do not know how the effect of “one gene regulating another” varies temporally. The number of time-lagged parameters in such a time series model can quickly grow with the maximum allowed time-lag (Table 1). In other words, because each of the gene expression varies over time, we need to know how the expression of one gene at one point in time regulates another gene’s expression at a future time in a parameter efficient manner. This issue is unique to Granger-Causality over time series data in scenarios with limited data, relative to the high dimensionality of permissible time-lagged covariate distributions.

A common way to address this issue has been to be conservative and train prediction models over large time windows to capture any long-term effects of covariates. This approach often runs into data sparsity issues and poor granger-causal link discovery accuracy (6). To deal with this issue, prior graphical granger methods (7) have artificially imposed sparsity constraints on the model parameters forcing co-efficients to collapse to zero, thereby reducing time-lagged parameters to estimate (from 600 – 700 to  $\leq 100$ ). However, Granger-Causality was not designed for large temporal networks with millions of time-lagged parameters; in fact, economists have often warned against blindly applying it over a large number of variables (8). One of the issues in applying Granger-causality directly to large temporal networks, is that large window sizes also result in the increase in chances of a violation of the positivity assumption (9, 10) - a condition necessary for consistency of the Graphical Granger methods, i.e certain time-lagged covariate

### Significance Statement

Granger-causal links discovered using observational time series are critical for many real-world applications in physical, biological and social sciences. The primary goal of this paper is to improve the discovery of Granger-causal links by learning to defer temporal model predictions through conditioning on covariate overlap using an uncertainty measure. By optimizing the prediction loss with a bounded threshold of overlap-based uncertainty over a subset of the data, we discover more granger-causal links critical for scientific hypotheses generation through accurate and generalizable temporal prediction models.

Authors A.B, S.J, and L.S contributed towards the conceptualization, investigation, methodology, writing and supervision of the paper. A.B contributed towards the implementation, visualization; and L.S contributed the funding acquisition for the work.

Prof. Subramanian is a co-founder of Entrupy Inc, Velai Inc, and Gaius Networks Inc and has served as a consultant for the World Bank and the Governance Lab. Prof. Jagabathula is a co-founder of Velai Inc. Velai Inc broadly works in the area of socio-economic predictive models. Mr. Balashankar is a Ph.D student at New York University, and is also funded in part, by the Google Student Research Advising Program. No other disclosures were reported.

\*Corresponding author e-mail: lakshmics.nyu.edu

distributions have been rarely or never observed previously and hence extrapolating the granger-causal links to such covariate distributions may be erroneous. In such scenarios, it might be best to defer prediction rather than predicting by extrapolating incorrectly.

To deal with the above challenges, in this paper, we propose a methodology to improve the recall of Granger-causal links by conditionally allowing for prediction deferrals. Specifically, given an outcome variable, a treatment variable, and a collection of covariates which are known to have passed the bivariate Granger-causal test (all of which are time series), our method parametrizes each of the Granger-causal links with two parameters: (a) the maximum window size,  $\delta$ , and (b) the variance threshold,  $\rho$ . The maximum window size specifies a bound on how long the treatment effect lasts and the variance threshold specifies when predictions are deferred. In particular, if the variance of our estimator for a given covariate value is above  $\rho$ , then we defer the prediction, suggesting that we do not have sufficient confidence in whether the treatment takes effect under the given covariate value based on the observations. Our method chooses the values of  $\delta$  and  $\rho$  to optimize a Granger-causal link recall metric, which is computed using only those covariate distributions whose variance is below the threshold  $\rho$ . A small value of  $\delta$  will result in fewer deferrals during training and thus a more reliable estimate of the recall metric, but might miss several links with longer time lags yielding lower recall, and the model becomes too sensitive to perturbations over smaller time windows. Very large values of  $\delta$  also result in lower recall because they result in a large number of deferrals during training and consequently less evidence from data to support link discovery. Thus our formulation trades-off the model's temporal sensitivity and overlap-based robustness, and learn predictive models that have high accuracy and are consistent with the known Granger-causal links.

Prior work does not consider link recovery but instead focuses on optimizing prediction accuracy using general purpose sparsity inducing techniques. In particular, multivariate Auto-Regression linear models (VAR) are trained to optimize prediction accuracy while inducing sparsity in the time lag parameters through Group Lasso penalty (4) regularization; links are included by comparing the prediction accuracy of the model with and without the treatment variable (11) to test for significance. The non-linear version of the above VAR Granger Causality models have also been proposed (12) which could model additive effects of the past of each series in a decoupled manner. Sequence prediction models (13) and graph attention models (14) which model the neighborhood of nodes to learn Granger-causal links have also been studied. We build off of these constraints and demonstrate that augmenting the condition of overlap violation ensures that prediction models which learn to defer when specific covariate distributions have not been previously observed are better at discovering Granger-causal links.

Learning such robust time-lagged Granger-causal models can be of immense importance in various real world scenarios of causal discovery. For example, in our running example of gene expression networks, using time series data from multiple such experiments carried out in laboratory settings along with our Granger-causal parametrization framework, we can extract optimal lag parameters between different genes to understand when (time-lag) and how (covariate overlap) they

influence each other's expression. Similarly, in the human motion capture MoCAP task (15), we are able to improve the area under the precision-recall curve (AUC-PR) of detecting Granger-causal links in the human activity recognition dataset than baseline Granger-causality (12, 16) methods. Finally, given time series of Granger-causal news events, we improve monthly forecasting accuracy of stock prices. Each of these three tasks have a large set of time-lagged parameters (113K - 2.9M as per Table 1), but aim to predict a very small number of target variables. In such scenarios, the problem of over-parameterization may lead to prediction models that spuriously rely on multiple treatment time series variables. We overcome this limitation of prior sparsity inducing methods, by directly optimizing for the recall of Granger-causal links with sufficient covariate overlap.

The conditional temporal prediction models we have developed are applicable to a diverse set of forecasting tasks. Specifically, our conditional covariate based training approach has

- Reduced the number of parameters to learn by 2-3 orders, and achieved 25% better AUC-PR in discovering Granger-causal links than comparative baselines
- Improved prediction accuracy over held-out time series by 18-25% across three datasets in MoCap activity recognition, DREAM3 gene regulatory networks detection and the New York Times news-based stock price prediction tasks.
- Formalized the trade-off between time-lag sensitivity and overlap-based robustness and showed that artificially increasing the maximum time-lag leads to an over-specified model with sub-optimal link recovery and prediction accuracy.

## 2. Granger Causal Link Recovery

**A. Problem Setup.** We consider a time-series forecasting problem where the goal is to predict the value  $y(t+1)$  at time  $t+1$  of an outcome times series using the past observations  $\mathbf{X}(t, \delta) = \{\mathbf{x}_1(t, \delta), \mathbf{x}_2(t, \delta), \dots, \mathbf{x}_n(t, \delta)\}$ , the treatment time series  $\mathbf{v}(t, \delta)$ , the historical outcomes  $\mathbf{y}(t, \delta)$ , each of them a vector of values evaluated up to  $\delta$  discrete time steps back in time from timestamp  $t$ . We assume we are given a predictive model  $m$ , which outputs the future values of the outcome time series from the past observations:

$$\hat{y}(t+1) = m(\mathbf{X}(t, \delta), \mathbf{v}(t, \delta), \mathbf{y}(t, \delta)).$$

Each of the above variables  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n, \mathbf{v}, \mathbf{y}\}(t, \delta)$  are time-lagged multivariate variables with  $\delta$  time-lagged values going back from time instant  $t$ . This means the predictive model has a total of  $(n+2) \cdot \delta$  variables as input to predict the value of the outcome at time  $t+1$ . We also assume we are given non-parametric links of interest of the form  $\mathbf{v} \rightarrow \mathbf{y}$ . We now train this predictive model on historical observational data and use it to make fine-grained predictions for each covariate value  $(\mathbf{x}, \mathbf{v})$ . Implementing this approach requires us to answer two questions:

1. How long does the treatment effect last?
2. Under what covariate values does the treatment take effect?

Answering the first question allows to set the correct value of  $\delta$ . The second question is important when dealing with observational data because the positivity assumption in covariate overlap  $P(\mathbf{v}(t, \delta) = v | \mathbf{X}(t, \delta) = \mathbf{x}) > 0$  is often not met and blindly assuming it holds can lead to incorrect extrapolation. Prior work has focused on answering the first question using sparsity inducing methods (7), while we argue that the second question is of equal importance to learn a robust time-lagged Granger-causal model and has implications on answering the first question. We formulate this as a joint learning problem with the *goal* to learn a compact Granger-causality aware time series predictive model which also learns to defer from making over-confident predictions for time periods with very few conditional treatment and covariate observations in the data (17).

**B. Time-Lagged Granger-Causal Model Assumptions.** Learning the time-lag parameters based on temporal predictions have been studied in the domain of Granger causality (16), where the links are established based on the lag parameter in a time series of the causal variable that provide the highest reduction in regression error of predicting the effect variable in a multivariate setting. These models make assumptions of sparsity, i.e for a given  $(\mathbf{v}, \mathbf{y})$  only a small number of time-lagged variables of  $\mathbf{X}, \mathbf{y}, \mathbf{v}$  are predictive of future values of  $y(t+1)$ . Many sparsity enforcing methods have been proposed like the Lasso regularization (7) which minimize the number of non-zero weights in a linear model (12, 18), or propose an auxiliary task based regularization of jointly predicting the causal graph and optimal predictors (19) or propose a recall-based regularization method to model autocorrelated time series with latent confounders (20). One *overlooked assumption* in the above approaches in overcoming the overparameterization issue is that of the positivity assumption in covariate overlap ( $P(\mathbf{v}(t, \delta) = v | \mathbf{X}(t, \delta) = \mathbf{x}) > 0$ ) (21), between treatment and outcome, given observational data. We overcome the limitation of prior methods by addressing violations of the positivity assumption explicitly through covariate conditional variance estimation.

**C. Covariate Conditional Variance Estimation.** Under the ignorability (no unobserved confounders) assumption, we can either estimate the importance of the treatment on the outcome variable from observed data by either slicing data based on treatment (22) or incorporating the treatment variable as another covariate (17, 23, 24). In Algorithm 1, we adopt the latter approach and learn a time-lagged prediction model that optimizes the recall of Granger-causal links with the target variable  $y$ , using the covariates  $\mathbf{X}$  and the treatment variable  $\mathbf{v}$ . Now we explain the Granger-causal link significance test used, and how we use a conditioned version of it to incorporate the covariate overlap assumption. Finally, we tie all of these components into a Bayesian optimization algorithm that maximizes the recall of the Granger-causal links.

**Covariate Conditional Treatment Importance:** To overcome the intractability of ensuring the overlap assumption for large number of covariates, we use the approach used by (17). We estimate the lack of overlap using the conditional variance -  $\hat{V}(\mathbf{x}, \delta)$  of the predictive models  $\hat{m}, \tilde{m}$ , which predict  $\hat{y}(t+1)$  (Eqn 5). We use a non-parametric Conover Squared Ranks (SR) test-statistic used for testing for equality of variance in prediction errors (11), to approximate the *new*

*information* that helps in improving the prediction accuracy of models. Once the prediction models are trained on certain splits of the data for a given outcome, we then estimate the variance by the bootstrapping method and evaluating the covariate variance on numerous held-out development splits. For a given threshold  $\rho$ , we adopt a trimming policy (25) with the rejection policy, conditioning on covariates  $\mathbf{x}$ , where the covariate variance  $\text{Var}[\hat{V}(\mathbf{x}, \delta)]$  is above a threshold  $\rho$  (17), and then compute the *Time-Lagged Conditional Treatment Importance - TCTI*( $\delta, \rho$ ).

$$TCTI(\delta, \rho) = \mathbb{E}_{\mathbf{x}: \text{Var}[\hat{V}(\mathbf{x}, \delta)] \leq \rho} \hat{V}(\mathbf{x}, \delta) \quad [1]$$

This formulation outlines that when the predictive causality test statistic has high variance, we should not rely on those slices as they violate the overlap assumption. If the estimated value of TCTI passes the SR test with  $\alpha$  significance ( $\alpha = 0.05$ ) under a  $F(1, (n+2)\delta+1)$  distribution, we then consider the treatment variable  $v$  to have an  $\alpha$ -significant Granger-causal link with  $y$ . This process is then repeated for all treatment variables to give the recall  $w(\delta, \rho)$ :

$$w(\delta, \rho) = \frac{\# \text{ links } \alpha\text{-significant with } (\delta, \rho) \text{ as per Eqn 1}}{\# \text{ Granger-causal links}} \quad [2]$$

For a given time-sensitivity lag parameter, the trade-off with overlap assumption to optimize recall can be understood by varying the hyper-parameters together. We implemented 3 approaches to fine-tune the hyper-parameters  $\delta, \rho$ , including a grid search, random search and a Bayesian optimization technique (26) outlined in Algorithm 1, that directly models the recall of Granger-causal links as a utility function  $w(\delta, \rho)$  (see detailed methods).

---

**Algorithm 1** BayesOpt for Exploring the Time-Overlap Trade-offs

---

- 1:  $(\delta, \rho) \sim U$ : uniform random distribution,  $\hat{\delta} \leftarrow 0$
  - 2:  $M = \{\}, \hat{w} \leftarrow 0, \psi \leftarrow \text{stopping criterion}$
  - 3: **while**  $\hat{w} < 1 - \psi$  **do**
  - 4:   Update BayesOpt acquisition function
  - 5:   Acquire  $D = \{k \text{ values of } \delta\}$  from BayesOpt
  - 6:   Update  $\max_{\delta} \leftarrow \max_{\delta \in D}(\delta, \max_{\delta})$
  - 7:   Train  $\hat{m}_{\max_{\delta}}, \tilde{m}_{\max_{\delta}}$  and update  $\hat{\delta} \leftarrow \max_{\delta}$  if  $\max_{\delta} > \hat{\delta}$
  - 8:   Acquire  $R = \{k \text{ values of } \rho\}$  from BayesOpt
  - 9:   Update  $f, \delta^*, \rho^*$  over  $R$  to maximize  $w(\delta, \rho)$ : Eqn (2)
  - 10: **end while**
  - 11: Return  $\delta^*, \rho^*$
- 

**Inference:** Once we obtain the optimal  $\delta^*, \rho^*$ , for each outcome variable, we use the trained models to make time-series predictions over unseen data. If the variance for a given covariate as pre-computed by the test statistic in Eqn 5 is above the threshold  $\rho^*$  at training time, we continue to defer to make predictions on those covariates at inference time. The resulting time-lagged prediction model then is used directly to infer the Granger-causal links based on the non-zero model parameters in the models (Sec C).



### 3. Results

We now demonstrate the efficiency of our approach in Granger-Causal link discovery, while showing that our prediction deferrals do not reduce prediction accuracy. Also, we show choosing the right set of temporal and covariate hyperparameters are critical for this improvement, as compared to generic sparsity inducing baselines where the choice is adhoc.

**A. Granger-Causal Link Discovery.** We demonstrate that our proposed method by optimizing  $TCTI(\delta, \rho)$  and hence the Granger-causal link recall, we improve both the recall of Granger-causal links and the prediction accuracy across 3 datasets and 4 baseline Granger causal models (see baselines Sec C in detailed methods). In Figure 1, we see that the prediction accuracy (y-axis) improves by 18-21% and the recall of Granger-causal links (x-axis) improves by 25% across 3 datasets for each of the four prediction models'  $TCTI\{-VAR, Neural\ Granger, Graph\ Generative, Graph\ Attention\}$ , when the time lag and overlap parameters are optimized with 25 random restarts. The baseline models base- $\{VAR, Neural\ Granger, Graph\ Generative, Graph\ Attention\}$  indicated as dots in the plot, have the sparsity constraint enforced through Lasso regularization loss, and end up with a lower recall of Granger-causal links i.e more Granger-causal links are not used for prediction, with low prediction accuracy. Further, to be comparable with baseline models which do not defer, we do not defer predictions at inference time, but only while learning the optimal time-lag and overlap parameters  $\delta, \rho$  in Algorithm 1.

**B. Importance of Prediction Deferrals.** To understand how deferring predictions on covariates with high variance as per  $TCTI(\delta, \rho)$  has helped learning better prediction models on the remaining covariates, we see that for each of the 3 tasks and 4 models in Figure 3, that there is an increase in prediction accuracy for the covariates we choose to predict as compared to the overall baseline model. Further, the increase is greater when the number of covariate slices deferred is greater. Thus by deferring predictions over covariates where overlap assumptions are violated, we improve the prediction accuracy and rely on robust Granger-causal links for prediction.

**C. Variation in Time-Lag and Overlap Parameters.** To understand how setting time-lag and overlap parameters is critical for the high performance of our approach, we plot the range of parameters required to achieve a fixed prediction accuracy across the 4 prediction models with our approach. In each of the dataset, we see that the distribution of parameters for time-lag, overlap-constraint for the links for a given Granger-causal model (the Graph Attention Model) has high variability as shown in Figure 4. Also, for a fixed prediction accuracy of outcomes, we see that the links for the 4 modeling choices are clustered into 3 groups - (low  $\delta$ , high  $\rho$ ), (low  $\delta$ , low  $\rho$ ), (high  $\delta$ , low  $\rho$ ). The lack of links with relatively high  $\delta$  and high  $\rho$  further empirically affirms the trade-off between time-lag and overlap-constraints as shown in Figure 5.

**D. Hyperparameter Optimization Method.** We also see that across the 3 datasets, and the 4 modeling techniques - the choice of the hyperparameter fine-tuning methodology can impact the number of Granger-causal links we can recover with significant  $TCTI(\delta, \rho)$  as shown in Figure 2. While grid search

and random sampling provide good initial estimates of the maximum recall that can be achieved, we see that BayesOpt quickly outperforms these brute force search mechanisms for the same number of model re-trainings. This drastically reduces the time and cost required to identify the temporal and overlap characteristics of the Granger-causal links.

### 4. Discussion

This paper provides a new framework for learning temporal and overlap parameters in Granger causal models for time series modeling tasks. These time-lagged models demonstrate significant gains compared to alternative formulations across three completely disparate time series tasks: news-based stock price prediction, the DREAM3 gene expression network analysis and MoCAP human motion recognition. The DREAM3 datasets and the associated inference research challenges have significant broader implications to the systems biology research community. Given our observed AUROC and AUC-PR gains, we hope that the systems biology community can benefit from our code base and model results (which we will release to the public). Similarly, the MoCAP dataset is widely used within the motion capture community. Here, we are able to demonstrate AUC-PR gains for detecting human activity and improve the recall on connecting these movements to Granger-causal links in the human skeletal structure. Our results on the stock market dataset would be highly relevant for the finance community and we believe this line of work can be extended to build causally-aware predictive models for socio-economic applications. Across all these open data sets, we believe our research conforms to the ethical guidelines outlined in these communities.

We have demonstrated the need to parameterize causal links with their associated temporal sensitivity with awareness of overlap assumption violations. In time series data, we show that there exists a trade-off between how temporally sensitive any prediction model incorporating the causal links can be while not compromising on the covariate overlap assumption. This allows us to further build better prediction models while not relying on data that lacks overlap while attempting to capture long term effects. Specifically, in the MoCap activity recognition task, we see that an inaccurate time-lag incorporated into the model can lead to inaccurate activity predicted which may have implications on applications in augmented and virtual reality. These errors emanate from incorrectly reconstructing the skeleton of the human body from the sensor data. In the DREAM3 gene regulatory network, if the expression of one gene is incorrectly predicted, then we misinterpret one gene interacts with another - which may lead to ineffective drug candidates that target the gene regulatory networks. Finally, in the stock price prediction task, an incorrect time-lag can be catastrophic for any algorithmic trading application that relies on news based indicators - such crashes in the stock market have been anecdotally reported as flash crashes. Hence in all these scenarios, using the correct time-lagged model has implications in downstream applications, and if left unaddressed can lead to spurious Granger-causal links being incorporated. Further, since the size of such temporal networks are quite large with millions of parameters of estimate, a principled way of addressing covariate uncertainty through prediction deferrals can further improve the trust in practitioners.

## 5. Materials and Methods

**A. Datasets. DREAM3** The DREAM3 gene expression network inference challenge (5) consists of 5 datasets, 2 for E.Coli and 3 for Yeast, with each dataset containing 100 time series. Each of the time series has 46 variables, each of them gene expression replicates observed at 21 time instants. For each of these 46 variables, we consider in a round-robin version that one of those variables is the outcome, one is the treatment, and the rest as covariates. This allows us to estimate  $TCTI(\delta, \rho)$  for a total of 2070 combinations of treatment and outcomes, given the covariates. The ground truth contains directed Granger-causal links between 46 replicates, and through our prediction models, we parametrize each of the Granger-causal links by sweeping over values of  $\delta, \rho$  with the maximum significant value of  $TCTI(\delta, \rho)$ . In parallel, we also report the AUROC (Area under the Receiver-Operator Curve) and AUC-PR (Area under the precision recall curve) for the classification task of detecting the binary gene expression.

**MoCAP** The CMU MoCAP dataset (15) consists of motion sensor data, for 54 joints, collected from two subjects for a total of 2024 time points. Here, we are given the Granger-causal links of the human skeleton and we learn the time lag and overlap parameters for the classification of each of the activities - jumping jacks, side twists, arm circles, etc. Here, we report the AUC-PR (Area under the Precision-Recall curve) for detecting the human activity based on the movements in the human joints, along with the recall of the Granger-causal links in the human skeleton.

**Stock Price Prediction** In the stock prediction task, the outcomes are each of the 10 stock prices from 2010-13 (with 2013 as the test split), and for the treatment variable, we are given the top 10 financial news based Granger-causal factors and a further 100 covariates extracted from NY Times by (18). Here too, we measure the Root Mean-Squared Errors (RMSE) of the predicted values and the fraction of time instants where we had to defer the prediction.

**B. Methods.** We now present the framework of our methodology to compute the trade-off parameters for the assumptions of temporal sensitivity and overlap-based robustness. By fine-tuning these parameters at development time to maximize recall of Granger-causal links, we learn the parameters that provide the best possible supporting evidence for the links in each of the domains. Note that sparsity inducing regularization constraints like Lasso, might lead to zero coefficient values for certain variables, and hence the recall of Granger-causal links, that measures whether all known Granger-causal links are used in the prediction model, may drop when optimizing for prediction accuracy.

### Time-Lagged Predictive Causality

Under the ignorability (no unobserved confounders) assumption, we can either estimate the importance of the treatment on the outcome variable from observed data by either slicing data based on treatment (22) or incorporating the treatment variable as another covariate (17, 23, 24). In this paper, we adopt the latter approach and learn a prediction model with high accuracy of the target variable  $y$ , using the covariates  $\mathbf{X}$  and the treatment variable  $\mathbf{v}$ . To overcome the intractability of ensuring the overlap assumption for large number of covariates, we use the approach used by (17) to estimate the lack of overlap using the conditional variance of the predictive model

$m$ , which predicts  $\hat{y}(t+1)$ . (Section C).

Since the values of importance weights can vary depending on the choice of models, for example in the case of linear regression - they are coefficients, whereas in non-linear network models, there are attention weights, activation vector alignment - similar to Granger methods (27), we use a non-parametric Conover Squared Ranks (SR) test for equality of variance (11), to test if the treatment variable provides any *new information* that helps with the prediction accuracy of models:  $\hat{m}$  with and  $\tilde{m}$  without the treatment variable as input. The test is run on the prediction errors  $\epsilon(\hat{y}_{(\mathbf{x}, \delta)}(t+1)), \epsilon(\tilde{y}_{(\mathbf{x}, \delta)}(t+1))$  produced by prediction models  $\hat{m}, \tilde{m}$  respectively on held-out temporally disjoint test data.

$$\hat{y}_{(\mathbf{x}, \delta)}(t+1) = \hat{m}(\mathbf{X}(t, \delta) = \mathbf{x}, \mathbf{v}(t, \delta), \mathbf{y}(t, \delta)) \quad [3]$$

$$\tilde{y}_{(\mathbf{x}, \delta)}(t+1) = \tilde{m}(\mathbf{X}(t, \delta) = \mathbf{x}, \mathbf{y}(t, \delta)) \quad [4]$$

$$\hat{V}(\mathbf{x}, \delta) = SR(\epsilon(\hat{y}_{(\mathbf{x}, \delta)}(t+1)), \epsilon(\tilde{y}_{(\mathbf{x}, \delta)}(t+1))) \quad [5]$$

The SR test we use, is the non-parametric alternative of the Levene's test (28), which itself is the robust alternative for non-normal distributions to the 1-way between-groups analysis of variance (ANOVA) (29) test to detect equality of population means. We use this test over the parametric ones as we do not make any assumption of the distribution of variance (normal), as our test of overlap violation cannot work if we already assume that there is an underlying normal distribution. The input for the prediction model are  $\delta$  lagged time series of covariates, treatments and outcomes, and we will control the length of this time series as part of our methodology. We vary the  $\delta$  and train jointly - warm starting hidden model parameters as the time lag  $\delta$  increases, instead of training a separate model from random initialization per value of  $\delta$ . Thus, we are able to compute the temporal lag that maximizes the prediction accuracy of the target variable that is Granger-causally linked to covariates. We characterize that a Granger-causal link to be supported by the observed data (or to be recalled), if adding a Granger-causal variable's temporal lag causes an increase in the prediction accuracy of the outcome conditioned on the covariate  $\mathbf{X}(t, \delta) = \mathbf{x}$ , as show by a test statistic with a p-value below the statistical significance threshold ( $\alpha = 0.05$ ) under a  $F(1, n+1)$  distribution, as compared to not incorporating the Granger-causal variable at all. Otherwise, we characterize that Granger-causal link as not yet observed in the data. To convert the time-sensitivity into a variance based estimate, we compare the prediction errors  $\epsilon$  and characterize it by the inequality of variance  $\hat{V}(\mathbf{x}, \delta)$  given by the test statistic of the Squared Ranks (SR) test. Models  $\hat{m}, \tilde{m}$  are trained and tested (Eqn 5) on temporally disjoint time series data.

### Overlap-based Conditional Treatment Importance

We now have to overcome the intractability of conditioning on the large number of covariates. Here too, we use the previously trained models:  $\hat{m}, \tilde{m}$  to predict the target variable, but use the variance of the treatment importance estimate (17). Once the prediction models  $\hat{m}, \tilde{m}$  are trained on certain splits of the data for a given outcome, we can then estimate the variance by a bootstrapping method and evaluating  $\hat{V}(\mathbf{x}, \delta)$  on numerous held-out development splits.

$$Var[\hat{V}(\mathbf{x}, \delta)] = Var[SR(\epsilon(\hat{y}_{(\mathbf{x}, \delta)}(t+1)), \epsilon(\tilde{y}_{(\mathbf{x}, \delta)}(t+1)))] \quad [6]$$

For a given threshold  $\rho$ , we adopt a trimming policy (25) with the rejection policy, conditioning on covariates where the variance of  $Var[\hat{V}(\mathbf{x}, \delta)]$  is above a threshold  $\rho$  (17), and then compute the *Time-Lagged Conditional Treatment Importance* -  $TCTI(\delta, \rho)$ . While the models trained are dependent on  $\delta$ , the computation of  $TCTI(\delta, \rho)$  is done after the training is completed, rather than at training time.

This formulation clearly outlines that when the predictive causality test statistic has high variance, we should not rely on those slices as they violate the overlap assumption. For a given time-sensitivity lag parameter, this trade-off with overlap assumption can be understood by varying the hyper-parameters together. The utility function  $w(\delta, \rho)$  is given by the fraction of the Granger-causal links in the given set that passes the difference in means t-test (with significance level  $\alpha$ ) that  $TCTI$  is different as compared to the null distribution. We see that for thresholds:  $\rho$ , such that the overlap condition  $\mathbf{x} : Var[\hat{V}(\mathbf{x}, \delta)] \leq \rho$  is satisfied for all covariates  $\mathbf{x}$ , then there would be no deferral, and such a  $TCTI(\delta, \rho)$  would directly evaluate the Granger-causality of the links, and hence  $w(\delta, \rho)$  would be equal to 1.

We now outline the 3 approaches we undertake to fine-tune the hyper-parameters  $\delta, \rho$ .

**Grid Search:** By using a grid search for values of  $\delta \in \{1, 2, \dots, T\}$ ,  $\rho \in \{\eta, 2\eta, \dots, k\eta\}$ , for each Granger-causal link in the dataset, we search for the value  $\delta^*, \rho^*$  that maximizes the  $w(\delta, \rho)$ . This way, we search among all Granger-causal links, the sensitivity and robustness parameters that is best supported by the observed data. This can be time-consuming to be done for each model and we can upper-bound the time lag parameter  $T$  to train the model.

**Random Search:** Instead of an exhaustive grid search, in this approach, we sample values of  $\delta, \rho$  from a uniform distribution and choose the parameters that maximizes the  $w(\delta, \rho)$ . Here, we were able to heuristically choose a bound larger than  $T, k\eta$  respectively and can search among values not explored by the grid search. Although we can control the number of hyper-parameters to train and evaluate the models against, the computational cost in training the models remain.

**Bayesian Optimization:** To learn which hyper-parameters  $\delta, \rho$  result in high recall of Granger-causal links as per the significance level  $\alpha = 0.05$ , we used Bayesian optimization with the probability prior  $f$  parameterized by  $\theta$  to be drawn from Gaussian processes. Specifically, we maximize the fraction of links  $w$  validated with a significance level for a given value of  $(\delta, \rho)$ . The covariance kernel chosen is the ARD Matern 5/2 Kernel (26), which has been demonstrated to capture realistic hyper-parameter distributions in neural networks, while resulting in sampling functions that are twice differentiable. We choose hyper-parameters in parallel using the Bayesian roll out method, where the acquisition function is optimized the utility of expected improvement per trial.

As noted in Section B,  $\delta$  requires a higher cost and time as it requires re-training of the model, while  $\rho$  can be fine-tuned post-training. These costs are modeled independent of the hyper-parameter distributions by calculating the expected inverse duration of computation and incorporating it in the expected improvement per second utility.

**Inference:** Once we obtain the optimal  $\delta^*, \rho^*$ , for each outcome variable, we use the trained models  $\hat{m}_{\delta^*}$  to make time-series predictions over unseen data. If the variance

$\hat{V}(\mathbf{x}, \delta^*) > \rho^*$  for a given covariate  $\mathbf{x}$  as pre-computed by the test statistic in Eqn 5 using  $\hat{m}_{\delta^*}$  at training time, we defer to make predictions on those covariates at inference time. Thus, we can now compare the predictions made by our overlap-aware Granger-causal model with baselines on the data slices where we do not defer.

**C. Baselines.** We would like to present a few comparable baseline prediction models built to incorporate Granger-causality for time series and discuss how compact time-lagged versions of these models have been learnt. These models form the baselines on which we evaluate in Section 3.

**Multivariate linear auto-regressive models:** Multivariate Auto-Regression linear models (VAR) take the time series of the treatment, covariates and the lagged values of the target variable as input to predict the target variable for future time instants. Here, we compare the prediction accuracy of the model with and without the treatment variable using the non-parametric Squared Rank test (11) to test the significance of a non-zero coefficient in matrix C. Additionally, we also add the Group Lasso penalty (4) that has been shown to overcome the need of precisely estimating the time lag by applying the Lasso regularization.

$$y(t+1) = A \cdot y(t, \delta) + B \cdot \mathbf{X}(t, \delta) + C \cdot \mathbf{v}(t, \delta) + D \quad [7]$$

**Neural Granger causality:** The non-linear version of the above VAR Granger Causality models have also been proposed (12) which could model additive effects of the past of each series in a decoupled manner. Here, by modeling the task of Granger causality using componentwise multilayer perceptrons and recurrent neural networks, all time series are captured in an input layer of the neural network having a total of  $\delta \cdot n^2 \cdot W$  parameters, where W is the number of hidden units in the input layer. In order to model long time lags in Granger causality, they use component-wise recurrent neural networks (RNN) for each time series. Similar to the linear model, to enforce sparsity, Lasso penalty and the hierarchical group Lasso penalty have been proposed, which chooses a suitable lag for each of the time series - but ignores the covariate overlap violation.

**Generative Graph Neural Networks:** Another approach proposed in (13), is to model this as a sequence prediction by reducing the graph to Breadth-First-Search (BFS) based deterministic sequence. They use a hierarchical graph RNN structure to first model the node prediction problem. In our case, although we know all the nodes of the network ahead of time, we can use the edge prediction model and predict edges in a BFS sequence. While this is comparable to the Neural Granger Causality model, the number of parameters to be learnt is lesser:  $\delta \cdot n \cdot W$ .

**Graph Attention Model:** Also recently, with the success of attention models in natural language tasks like machine translation, attending over the neighborhood of nodes, instead of recurrent architectures has been shown to be specifically relevant for graphical causal modeling (14). This approach requires only the neighborhood of nodes and scales better than spectral representations of the graph, which need to be aware of the entire graph structure.

**Other Related Work:** We aim to understand the assumptions required to identify the optimal time sensitivity



parameters of Granger-causal links in time series data, once the direction and presence of the Granger-causal links have already been defined. Prior work in this space has focused on methodologies to increase recall of the causal links in auto-correlated time series (30) or regularize over unseen parts of the causal graph (19). However, such methods fail to quantify when it might be even possible to recover the optimal time lag parameters in an observed data distribution. The covariate relationship in time series have been explored in Granger causality with group boosting methods (4) or Markov random field regression (31) which capture the non-linear group information between variables in a time series. Prior methods (32) that maintain a set of probabilistic causal models and perform model selection can also benefit from the quantification of the trade-off between overlap and temporal parameters in longitudinal data.

**ACKNOWLEDGMENTS.** We thank Prof. Joan Bruna for his comments on a preliminary version of the manuscript.

## Data Sharing Plans

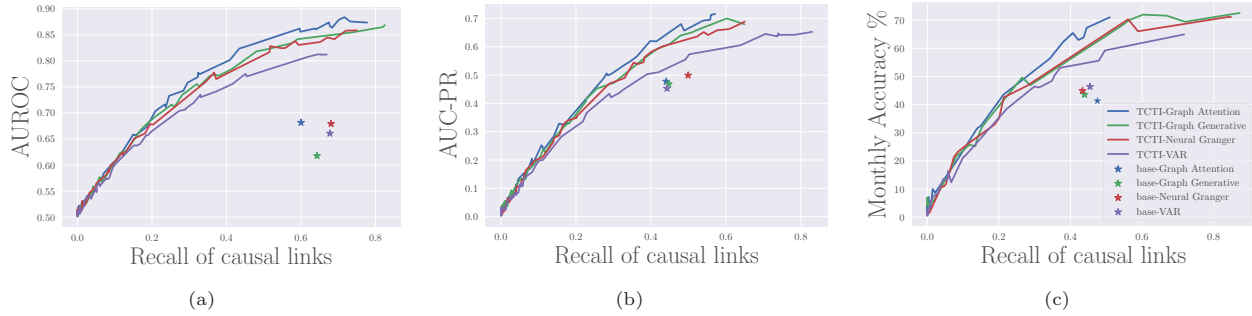
The data used for training and evaluating our methodology is publicly available, and the code used for generating the analysis will be made public.

1. CW Granger, Investigating causal relations by econometric models and cross-spectral methods. *Econom. journal Econom. Soc.* pp. 424–438 (1969).
2. A Arnold, Y Liu, N Abe, Temporal causal modeling with graphical granger methods in *Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '07. (ACM, New York, NY, USA), pp. 66–75 (2007).
3. AC Lozano, et al., Spatial-temporal causal modeling for climate change attribution in *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*. pp. 587–596 (2009).
4. AC Lozano, N Abe, Y Liu, 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*, KDD '09. (Association for Computing Machinery, New York, NY, USA), p. 577–586 (2009).
5. RJ Prill, et al., Towards a rigorous assessment of systems biology models: The dream3 challenges. *PLOS ONE* **5**, 1–18 (2010).
6. P Spirtes, An anytime algorithm for causal inference. in *AISTATS*. (2001).
7. A Arnold, Y Liu, N Abe, Temporal causal modeling with graphical granger methods in *Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '07. (Association for Computing Machinery, New York, NY, USA), p. 66–75 (2007).
8. J Peters, D Janzing, B Schölkopf, Causal inference on time series using restricted structural equation models in *Advances in Neural Information Processing Systems*. pp. 154–162 (2013).
9. PR Rosenbaum, DB Rubin, The central role of the propensity score in observational studies for causal effects. *Biometrika* **70**, 41–55 (1983).
10. VN Vapnik, *The Nature of Statistical Learning Theory*. (Springer), Second edition, (1999).
11. W Conover, R Iman, [rank transformations as a bridge between parametric and nonparametric statistics]: Rejoinder. *Am. Stat. - AMER STATIST* **35**, 124–129 (1981).
12. A Tank, I Covert, N Foti, A Shojaie, E Fox, Neural Granger Causality for Nonlinear Time Series. *ArXiv e-prints* (2018).
13. J You, R Ying, X Ren, WL Hamilton, J Leskovec, Graphrnn: A deep generative model for graphs. *CoRR abs/1802.08773* (2018).
14. P Veličković, et al., Graph attention networks in *International Conference on Learning Representations*. (2018).
15. CMU, Carnegie mellon university motion capture database (2009).
16. CWJ Granger, *Essays in econometrics*, eds. E Ghysels, NR Swanson, MW Watson. (Harvard University Press, Cambridge, MA, USA), pp. 31–47 (2001).
17. A Jesson, S Minderhmann, U Shalit, Y Gal, Identifying causal-effect inference failure with uncertainty-aware models in *Advances in Neural Information Processing Systems*, eds. H Larochelle, M Ranzato, R Hadsell, MF Balcan, H Lin. (Curran Associates, Inc.), Vol. 33, pp. 11637–11649 (2020).
18. A Balashankar, S Chakraborty, S Fraiberger, L Subramanian, Identifying predictive causal factors from news streams in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. (Association for Computational Linguistics, Hong Kong, China), pp. 2338–2348 (2019).
19. T Kyono, Y Zhang, M van der Schaar, Castle: Regularization via auxiliary causal graph discovery in *Advances in Neural Information Processing Systems*, eds. H Larochelle, M Ranzato, R Hadsell, MF Balcan, H Lin. (Curran Associates, Inc.), Vol. 33, pp. 1501–1512 (2020).
20. JL Gamella, C Heinze-Deml, Active invariant causal prediction: Experiment selection through stability in *Advances in Neural Information Processing Systems*, eds. H Larochelle, M Ranzato, R Hadsell, MF Balcan, H Lin. (Curran Associates, Inc.), Vol. 33, pp. 15464–15475 (2020).

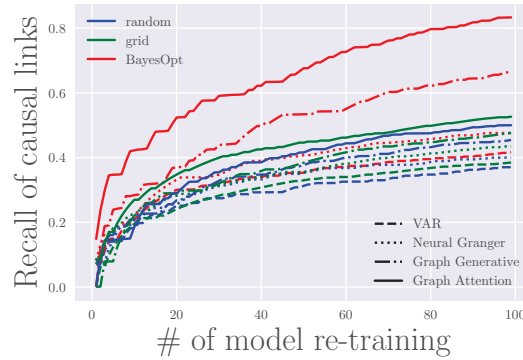
21. GW Imbens, Nonparametric estimation of average treatment effects under exogeneity: A review. *Rev. Econ. statistics* **86**, 4–29 (2004).
22. U Shalit, FD Johansson, D Sontag, Estimating individual treatment effect: generalization bounds and algorithms (2017).
23. A Gelman, J Hill, *Causal inference using regression on the treatment variable*, Analytical Methods for Social Research. (Cambridge University Press), p. 167–198 (2006).
24. CM Lawes, et al., Statin use in copd patients is associated with a reduction in mortality: a national cohort study. *Prim. Care Respir. J.* **21**, 35–40 (2012).
25. MA Hernán, JM Robins, Causal inference (2010).
26. J Snoek, H Larochelle, RP Adams, Practical bayesian optimization of machine learning algorithms in *Advances in neural information processing systems*. pp. 2951–2959 (2012).
27. CW Granger, BN Huangb, CW Yang, A bivariate causality between stock prices and exchange rates: evidence from recent asianflu. *The Q. Rev. Econ. Finance* **40**, 337–354 (2000).
28. MB Brown, AB Forsythe, Robust tests for the equality of variances. *J. Am. Stat. Assoc.* **69**, 364–367 (1974).
29. RA Fisher, Statistical methods for research workers in *Breakthroughs in statistics*. (Springer), pp. 66–70 (1992).
30. A Gerhardus, J Runge, High-recall causal discovery for autocorrelated time series with latent confounders in *Advances in Neural Information Processing Systems*, eds. H Larochelle, M Ranzato, R Hadsell, MF Balcan, H Lin. (Curran Associates, Inc.), Vol. 33, pp. 12615–12625 (2020).
31. Y Liu, A Niculescu-Mizil, A Lozano, Y Lu, Learning temporal causal graphs for relational time-series analysis in *Proceedings of the 27th International Conference on International Conference on Machine Learning*, ICML'10. (Omnipress, Madison, WI, USA), p. 687–694 (2010).
32. A Riva, R Bellazzi, Learning temporal probabilistic causal models from longitudinal data. *Artif. Intell. Medicine* **8**, 217–234 (1996) Temporal Reasoning in Medicine.

Dataset	n	$\delta$	# Variables	# Targets
MoCAP	54	2,024	113,344	12
DREAM3	46	10,500	504,000	5
Stocks	100	29,200	2,978,400	10

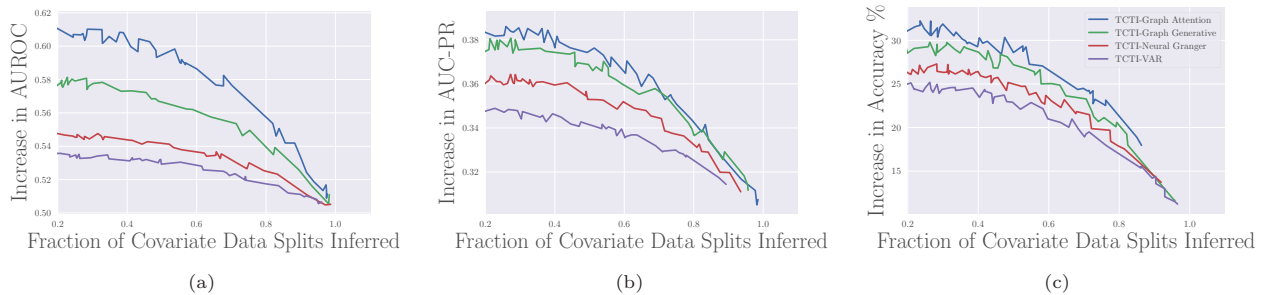
**Table 1. Problem of Over-parameterization in Time-Series Granger-Causal Models**



**Fig. 1. Recall of Granger-Causal Links vs Prediction Accuracy** (a) **DREAM3**: Across each of the 5 outcomes in the DREAM3 dataset, we see that as the recall of the Granger-causal links in the gene expression network increases, the AUROC of the time series of the gene expression level also increases. (b) **MoCAP**: As the recall of the Granger-causal links of the human skeleton network increases, the AUC-PR in the human activity recognition task increases. (c) **Stock**: As the recall of the Granger-causal links of the financial news factors increases, the prediction accuracy of 10 stock prices increases.

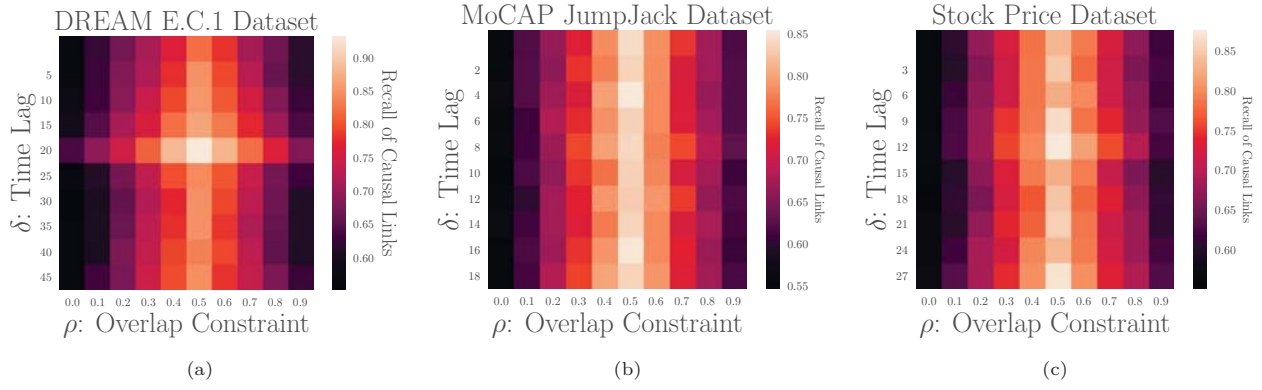


**Fig. 2. Sampling efficiency among hyperparameter search methods** The number of re-trainings required to fine-tune the time sensitivity and robustness overlap parameters to improve the recall of Granger-causal links in the DREAM gene expression dataset

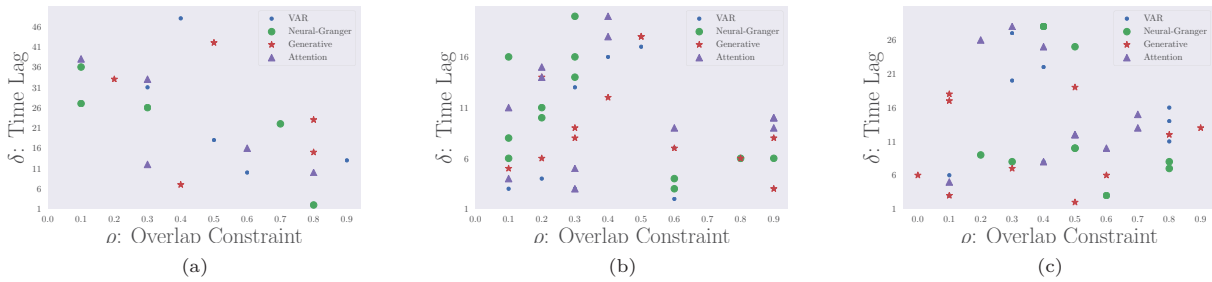


**Fig. 3. Prediction Deferrals effect on Accuracy** Choosing  $\delta$  based on an overlap based constraint, the prediction accuracy increases on the remaining test samples on (a) DREAM3 (b) MoCAP (c) Stock datasets for 4 Granger causal models. As the prediction models choose to defer on larger fractions of the covariate data splits, the increase in accuracy is higher.





**Fig. 4. Variation between time lag and overlap parameters for the Graph Attention Model on 3 datasets shows the need to learn them jointly**



**Fig. 5. Trade-off for a fixed prediction accuracy** Time sensitivity and robustness overlap among the fine-tuned Neural Granger Causal prediction models across Granger-causal links that provide the highest  $TCTI(\delta, \rho)$  for (a) DREAM3, (b) MoCAP and (c) Stock Price prediction datasets