

Conditional Independence Testing in the Presence of Temporal Correlation and Nonstationarity

Michael Wieck-Sosa* Michel F. C. Haddad† Aaditya Ramdas‡

March 7, 2025

Abstract

Understanding the structure of evolving temporal systems is a key goal in many scientific and engineering disciplines. While the standard toolkit for multivariate time series analysis — primarily based on linear vector autoregressive models — has many advantages, it can be difficult to capture complex nonlinear dynamics with these tools. This has motivated researchers to develop methods for variable selection, causal discovery, and graphical modeling for nonlinear time series. These methods are most effective when conditional dependencies are identified accurately, which is why tests for conditional independence are frequently used. Unfortunately, many nonparametric tests for conditional independence rely on techniques like permutation and sample splitting, making it unclear how they can be adapted to the nonstationary time series setting. The few tests that do accommodate temporal dependence require stationarity, which may be an unrealistic assumption to make even after differencing, detrending, or transforming the time series to growth rates.

In this paper, we introduce a general framework for conditional independence testing that is robust to both nonstationarity and temporal dependence. As far as we know, this is the first framework to enable conditional independence testing with a single realization of a nonstationary nonlinear process. The key technical ingredients are time-varying regression estimation, time-varying covariance estimation, and a distribution-uniform strong Gaussian approximation for nonstationary nonlinear processes.

Contents

1	Introduction	2
2	The Dynamic Generalized Covariance Measure (dGCM)	5
2.1	Setting and notation	5
2.2	The null hypothesis of conditional independence	6
2.3	Time-varying regression functions	8
2.4	Main ideas of our work and the algorithm	8
3	Assumptions and the Theoretical Result for dGCM	11
3.1	Nonstationary observed processes	11
3.2	Prediction processes	13
3.3	Nonstationary error processes	14
3.4	Assumptions on dependence and nonstationarity	15
3.5	Theoretical result for dGCM	16

*Department of Statistics and Data Science, Carnegie Mellon University. Email: mwiecksosa@cmu.edu.

†Department of Business Analytics and Applied Economics, Queen Mary University of London.

‡Department of Statistics and Data Science, Machine Learning Department, Carnegie Mellon University.

4	dGCM with Sieve Time-Varying Regression (Sieve-dGCM)	18
4.1	Setting and notation	18
4.2	Locally stationary observed processes	19
4.3	Sieve time-varying nonlinear regression estimator	20
4.4	Locally stationary error processes	23
4.5	Assumptions on dependence and nonstationarity	24
4.6	Assumptions on local long-run covariances	27
4.7	Theoretical result for Sieve-dGCM	28
5	Numerical Simulations	29
5.1	Parameter selection via subsampling and minimum volatility	29
5.2	Analysis of level and power	31
5.3	Discussion	39
5.4	Accuracy of the strong Gaussian approximation	39
6	Future Work	41
A	Extensions	50
A.1	Alternative test statistics	50
A.2	Cyclostationary processes	50
A.3	Simulation-and-regression for nonstationary processes	50
A.4	Simplifications under stationarity	51
A.5	Additional tests for locally stationary processes	52
A.6	Piecewise locally stationary processes	54
A.7	Weakening the assumptions on the error processes	54
B	Distribution-Uniform Theory	54
B.1	Literature review of distribution-uniform inference	54
B.2	Distribution-uniform strong Gaussian approximation	55
B.3	Distribution-uniform feasible Gaussian approximation	58
B.4	Auxiliary Lemmas	59
C	Proofs of Theoretical Results for dGCM and Sieve-dGCM	61
C.1	Proof of Theorem 3.1	61
C.2	Proof of Theorem 4.1	72

1 Introduction

A great deal of work has been dedicated to developing tests for conditional independence. That is, testing whether two random vectors X and Y are independent given a third random vector Z . For example, there are conditional independence tests based on conditional densities [157], characteristic functions [156], empirical likelihood ratios [158], discretization [106, 75], permutation [50, 145], kernels [57, 182, 146], copulas [21], and conditional mutual information [138]. Also, there are many conditional independence tests based on regressing X on Z and Y on Z followed by testing for independence between the residuals [118, 122, 125, 52, 180, 181].

Unfortunately, conditional independence tests oftentimes struggle to control the Type-I error in finite samples, as shown by Shah and Peters [149]. In fact, Shah and Peters [149] prove that conditional independence testing is *fundamentally impossible* without making further assumptions. This issue has sparked significant interest in conditional independence testing over the last several years. We begin by providing an overview of recent advances in conditional independence testing. Afterwards, we discuss how our work addresses limitations in the existing literature. Finally, we motivate our work by reviewing key applications of conditional independence tests for time series in areas such as variable selection, causal discovery, and Granger causality.

The hardness of conditional independence testing. The no-free-lunch result from Shah and Peters [149] states that if one wants to have a conditional independence test with Type-I error control for all absolutely continuous (with respect to the Lebesgue measure) triplets of random vectors (X, Y, Z) , then this conditional independence test cannot have power against any alternative hypothesis. To make the conditional independence testing problem feasible, we must consider a smaller subset of the null hypothesis and use domain knowledge to select an appropriate conditional independence test. This hardness result was revisited by Neykov, Balakrishnan, and Wasserman [115] and Kim et al. [85], and was extended to the time series setting by Bodik and Pasche [19].

Shah and Peters [149] proposed a conditional independence test based on the *generalized covariance measure* (GCM), which is a suitably normalized sum of the products of the residuals from the regressions of X on Z and Y on Z . In this case, the practitioner’s domain knowledge is used to select appropriate regression methods for the problem at hand. In contrast to the conditional independence tests previously mentioned, Shah and Peters [149] demonstrate that the GCM test has asymptotic Type-I error control, *uniformly* over a large collection of distributions for which the null hypothesis of conditional independence holds.

Since then, numerous tests have been developed which draw inspiration from the original GCM test [144, 100, 33, 164, 78, 25, 101]. Our conditional independence test can be considered a GCM-type test for the nonstationary nonlinear time series setting. As we will discuss, moving to this setting introduces several complexities and requires completely different theoretical tools than the original GCM test. Although we develop a GCM-type conditional independence test in this work, we point the reader to another influential literature about conditional independence testing based on versions of the model-X assumption [24, 96, 116, 8, 73, 9, 147, 65].

Limitations of the existing literature. Most of the previously discussed conditional independence tests lack Type-I error control guarantees outside the iid setting. Furthermore, the literature on conditional independence testing when given only a single realization of a nonstationary process remains strikingly limited. To the best of our knowledge, only two conditional independence tests have been proposed for this setting.

First, Malinsky and Spirtes [104] introduce a conditional independence test for nonstationary linear vector autoregressions with iid Gaussian errors. Specifically, they study processes that exhibit “stochastic trends” so that the first difference of the process is stable. In contrast, our conditional independence test allows for nonlinear processes with very general forms of nonstationarity and time-varying regression functions with non-iid and non-Gaussian errors. Moreover, we demonstrate that our conditional independence test possesses uniformly asymptotic Type-I error control, as established for the GCM test from Shah and Peters [149].

Second, Flaxman, Neill, and Smola [54] develop a conditional independence testing framework for non-iid data based on Gaussian process regression. The main idea is to pre-whiten the non-iid data using Gaussian process regression to control for dependencies (e.g. spatial, temporal, or network), which should yield iid residuals. The next step is to test for independence between these residuals using the Hilbert-Schmidt Independence Criterion (HSIC) [64]. The authors mention that their framework could be used with nonstationary covariance functions, although this idea was not developed further.

We also mention some conditional independence tests designed for the setting in which multiple realizations of a stochastic process are available. Manten et al. [105] develop a conditional independence test for stochastic processes using the signature kernel. Christgau, Petersen, and Hansen [33] introduced a framework for testing so-called “conditional local independence” relationships among point processes. Lundborg, Shah, and Peters [100] introduce a conditional independence test for function-valued random variables. Also, we note that there is a growing literature on (unconditional) independence testing for nonstationary processes. Liu et al. [98] develop independence tests based on the HSIC [64]. These tests require multiple realizations of the nonstationary process, whereas the independence tests for locally stationary processes from [23, 13] only require one realization.

Variable selection. A central problem in statistics and machine learning is variable selection. Conditional independence tests can be used for variable selection when paired with multiple testing procedures to control the false discovery rate (FDR) [15, 16]. In the context of forecasting, the goal is to identify a minimal subset $S \subseteq \{1, \dots, p\}$ out of p signals (including relevant lags) such that, for all times t , the forecasting target Y_{t+h} at horizon h is conditionally independent of the other signals $(X_t^i)_{i \notin S}$.

given $(X_t^i)_{i \in S}$. This minimal subset S is called a Markov blanket for Y_{t+h} ; see Pearl [120] and Candès et al. [24] for more discussion. We contribute to this literature by providing a conditional independence test flexible enough to be used for identifying relevant forecasting signals in unstable environments.

Causal discovery. The discovery of time-lagged causal relationships (see Figure 1) from observational time series is an important problem in numerous scientific domains [141]. Conditional independence tests for time series are a core component of constraint-based and hybrid causal discovery algorithms designed for temporally correlated data. For example, Runge et al. [140] used the conditional mutual information-based conditional independence test from Runge [138] in a causal discovery algorithm for time series called PCMCI, which builds on the classical PC algorithm [153]. Over the last several years, causal discovery for *nonstationary* time series has become an increasingly active area of research [104, 72, 53, 49, 142]. We stress that the conditional independence test used in the causal discovery algorithm must be appropriately tailored to the characteristics of the data. If the underlying conditional independence test fails to account for temporal dependence or nonstationarity, then the causal discovery algorithm may produce incorrect conclusions about the causal structure of the process. Our work fills a relevant gap in the literature on causal discovery for nonstationary nonlinear time series by providing a practical conditional independence test for this setting.

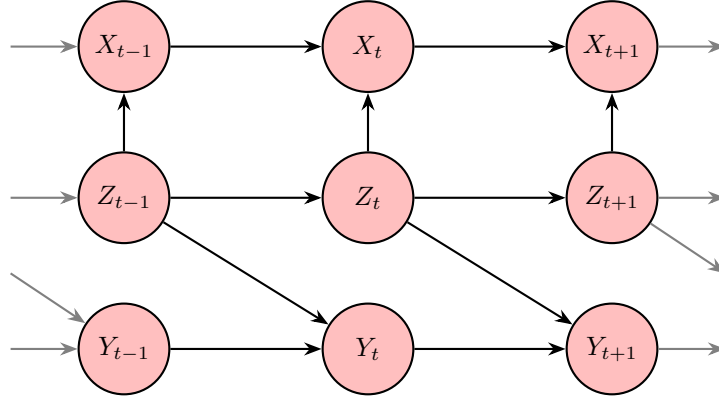


Figure 1: Causal graph depicting the time-lagged causal relationships among the stochastic processes $X = (X_t)_{t \in \mathbb{Z}}$, $Y = (Y_t)_{t \in \mathbb{Z}}$, and $Z = (Z_t)_{t \in \mathbb{Z}}$. The causal graph shows that Z is a common cause of both X and Y , directly influencing X in the same time period and affecting Y with a one time step delay. In this example, the causal structure of the multivariate process (X, Y, Z) remains fixed over time, though the causal effects themselves may vary over time.

Granger causality. Lastly, we discuss a commonly used framework for assessing relationships between stochastic processes called Granger causality. Unfortunately, the numerous definitions of Granger causality have caused a great deal of confusion. Recently, Shojaie and Fox [152] have written a comprehensive review to help clarify matters. We provide a highly condensed summary of Sections 2 and 4 in their review.

The original definition of Granger causality from Granger [62] is about prediction. Informally, a process X is said to be Granger causal of another process Y if the variance of the *optimal prediction* of Y_t at time t using all the relevant information up to time $t - 1$ is increased by removing the history of X up to time $t - 1$. See Section 2 in Shojaie and Fox [152] for the exact definition and the stringent conditions under which this predictive definition corresponds to genuine causality as in Pearl [119]. While this original definition does not assume linear dynamics, much of the following methodology revolves around the identification of coefficients in linear vector autoregressive (VAR) models with p time series [63, 103, 12].

Another definition of Granger causality, referred to as *strong* Granger causality [55], is stated in terms of conditional independence relationships among stochastic processes. Let $(X^i)_{i \in [p]}$ be p signals used to predict the target Y . The process X^i is said to be (strongly) Granger noncausal of Y if, for all times t , Y_t is conditionally independent of the history of the signal X^i up to time $t - 1$ given the history of the other signals $(X^j)_{j \in [p] \setminus \{i\}}$ up to time $t - 1$. See Definition 2 in Shojaie and Fox [152] for the exact definition, and the rest of Section 4 therein for more discussion.

Notably, Eichler [51] introduced a comprehensive graphical modeling framework for time series based on strong Granger causality, which can be detected using conditional independence tests for nonlinear time series [146, 21]. In a similar vein, our proposed conditional independence test can be used to detect strong Granger causality for nonlinear time series with *nonstationary* dynamics. This can be incorporated into graphical modeling frameworks for nonstationary nonlinear time series, analogous to Basu and Rao [11].

We note that there are also various techniques for assessing nonlinear Granger causality that do not utilize conditional independence testing. For instance, the neural Granger causality method from Tank et al. [159] extracts Granger causal structures by using sparsity-inducing penalties on the weights of structured multilayer perceptrons (MLPs) and recurrent neural networks (RNNs). Additionally, there is an influential strand of literature connecting Granger causality and directed information theory [3, 124]. See Section 4 of Shojaie and Fox [152] for more discussion of nonlinear Granger causality.

Paper outline. The rest of the paper is structured as follows. In Section 2, we discuss the main ideas and implementation of our proposed conditional independence test. In Section 3, we introduce the details of our theoretical framework. In Section 4, we develop the theoretical justifications for our test in the context of the sieve time-varying nonlinear regression estimator from Ding and Zhou [48] within the well-studied framework of locally stationary processes. In Section 5, we demonstrate the satisfactory performance of our test through comprehensive numerical simulations. Along the way, in Subsection 5.1, we introduce a novel cross-validation procedure based on subsampling for selecting the parameters of “global” estimators of time-varying regression functions. In Section 6, we discuss promising avenues for future work. In Section A, we consider several extensions of our conditional independence testing framework. In Section B, we state a distribution-uniform strong Gaussian approximation for high-dimensional nonstationary nonlinear processes.

2 The Dynamic Generalized Covariance Measure (dGCM)

In this section, we give a high-level overview of our work. Specifically, we introduce the notation, main ideas, and implementation of our proposed *dynamic generalized covariance measure* (dGCM) test. For expository purposes, we delay the technical details of our theoretical framework until Section 3.

2.1 Setting and notation

We work in a triangular array framework for high-dimensional nonstationary nonlinear time series. Let $(X_{t,n}, Y_{t,n}, Z_{t,n})_{t \in [n]}$ be the observed sequence of length $n \in \mathbb{N}$, where $[n] = \{1, \dots, n\}$. We use the notation $X_n = (X_{t,n})_{t \in [n]}$, $Y_n = (Y_{t,n})_{t \in [n]}$, $Z_n = (Z_{t,n})_{t \in [n]}$ to refer to sequences of length n , and we use the notation X, Y, Z to generically refer to the processes with any length. Let $d_X = d_{X,n}$, $d_Y = d_{Y,n}$, $d_Z = d_{Z,n}$ denote the dimensions, which can grow with n . Denote dimension $i \in [d_X]$ of $X_{t,n}$ by $X_{t,n,i}$, dimension $j \in [d_Y]$ of $Y_{t,n}$ by $Y_{t,n,j}$, and dimension $k \in [d_Z]$ of $Z_{t,n}$ by $Z_{t,n,k}$.

Next, we introduce notation for the time-offsets of each dimension of $X_{t,n}$, $Y_{t,n}$, $Z_{t,n}$ because we want to infer time-lagged conditional dependencies. Negative time-offsets are called *lags* of the process, and positive time-offsets are called *leads* of the process. Time-offsets of zero are allowed so that contemporaneous conditional dependencies can be considered. Let

$$A_i \subset \{-n+1, \dots, n-1\}, B_j \subset \{-n+1, \dots, n-1\}, C_k \subset \{-n+1, \dots, 0\},$$

be the sets of time-offsets of $X_{t,n,i}$, $Y_{t,n,j}$, $Z_{t,n,k}$ under consideration. We require the time-offsets C_k to be non-positive so that the conditioning variables are known at time t for purposes that will be made clear in the next few subsections. In practice, the largest (in magnitude) time-offsets should be selected small enough so that there is a sufficient amount of data to conduct the test.

Denote the time-offset $a \in A_i$ of $X_{t,n,i}$ by $X_{t,n,i,a} = X_{t+a,n,i}$, the time-offset $b \in B_j$ of $Y_{t,n,j}$ by $Y_{t,n,j,b} = Y_{t+b,n,j}$, and the time-offset $c \in C_k$ of $Z_{t,n,k}$ by $Z_{t,n,k,c} = Z_{t+c,n,k}$. Denote the sets of all time-offsets by $A = \bigcup_{i=1}^{d_X} A_i$, $B = \bigcup_{j=1}^{d_Y} B_j$, $C = \bigcup_{k=1}^{d_Z} C_k$, and largest (signed) time-offsets by $a_{\max} = \max(A)$, $b_{\max} = \max(B)$, $c_{\max} = \max(C)$, and the smallest (signed) time-offsets by $a_{\min} = \min(A)$, $b_{\min} = \min(B)$, $c_{\min} = \min(C)$.

Since we are interested in time-lagged conditional independence relationships, it is often useful to refer to the subset of original times,

$$\mathcal{T}_n = \{1 - \min(a_{\min}, b_{\min}, c_{\min}), n - \max(a_{\max}, b_{\max}, c_{\max})\} \subseteq \{1, \dots, n\},$$

in which *all* time-offsets of each dimension of $X_{t,n}$, $Y_{t,n}$, $Z_{t,n}$ are actually observed. Going forward, we will write $t \in \mathcal{T}_n$ instead of $t \in [n]$ because we are only using the subset of times in which all time-offsets are observed. Denote the first time of \mathcal{T}_n by $\mathbb{T}_n^- = \min(\mathcal{T}_n)$, the last time of \mathcal{T}_n by $\mathbb{T}_n^+ = \max(\mathcal{T}_n)$, and the cardinality of \mathcal{T}_n by $T_n = |\mathcal{T}_n|$. Note that if no negative time-offsets (i.e. lags) are used then $\min(a_{\min}, b_{\min}, c_{\min}) = 0$, and if no positive time-offsets (i.e. leads) are used then $\max(a_{\max}, b_{\max}, c_{\max}) = 0$. Hence, if only time-offsets of zero are used, then $\mathcal{T}_n = [n]$.

For all $t \in \mathcal{T}_n$, denote the vectors with all dimensions and time-offsets of interest by

$$\mathbf{X}_{t,n} = (X_{t,n,i,a})_{i \in [d_X], a \in A_i}, \quad \mathbf{Y}_{t,n} = (Y_{t,n,j,b})_{j \in [d_Y], b \in B_j}, \quad \mathbf{Z}_{t,n} = (Z_{t,n,k,c})_{k \in [d_Z], c \in C_k}.$$

Denote the dimensions of $\mathbf{X}_{t,n}$, $\mathbf{Y}_{t,n}$, $\mathbf{Z}_{t,n}$ by $\mathbf{d}_X = \sum_{i=1}^{d_X} |A_i|$, $\mathbf{d}_Y = \sum_{j=1}^{d_Y} |B_j|$, $\mathbf{d}_Z = \sum_{k=1}^{d_Z} |C_k|$, respectively. Also, denote the entire processes by

$$\mathbf{X}_n = (\mathbf{X}_{t,n})_{t \in \mathcal{T}_n}, \quad \mathbf{Y}_n = (\mathbf{Y}_{t,n})_{t \in \mathcal{T}_n}, \quad \mathbf{Z}_n = (\mathbf{Z}_{t,n})_{t \in \mathcal{T}_n}.$$

We allow the number of time-offsets to grow with n , that is, $A_i = A_{i,n}$, $B_j = B_{j,n}$, $C_k = C_{k,n}$ and $A = A_n$, $B = B_n$, $C = C_n$. However, we require that the largest (in magnitude) time-offset grows at a slower rate than n such that as $n \rightarrow \infty$ we have $\min(a_{\min}, b_{\min}, c_{\min})/n \rightarrow 0$ and $\max(a_{\max}, b_{\max}, c_{\max})/n \rightarrow 0$ so that the number of observed times $T_n \rightarrow \infty$ arbitrarily slowly.

Since we allow *both* the number of time-offsets and the number of dimensions to grow with n , we introduce the index set

$$\mathcal{D}_n \subseteq \{(i, j, a, b) : i \in [d_X], j \in [d_Y], a \in A_i, b \in B_j\},$$

which contains all of the indices for the dimensions and time-offsets of interest. Note that \mathcal{D}_n is specified by the user, and need not contain all possible combinations. Going forward, we will often refer to the dimension/time-offset tuple by $m = (i, j, a, b) \in \mathcal{D}_n$ to lighten the notation. The index set \mathcal{D}_n depends on the sample size n through the dimensions *and* the time-offsets, so its cardinality $D_n = |\mathcal{D}_n|$ may grow with n . Note that D_n reflects the *intrinsic dimensionality* of the problem and will appear frequently in the rest of the paper. In the “best case” scenario, we allow $D_n = O(T_n^{\frac{1}{\delta}})$. See (18) and the rest of Subsection B.2 for the full details about how quickly D_n can grow.

For each $n \in \mathbb{N}$, let \mathcal{P}_n be a collection of distributions for the processes, which we allow to change with n . For expository purposes, we delay the technical details about \mathcal{P}_n until the end of Subsection 3.1. Next, we state the null hypothesis.

2.2 The null hypothesis of conditional independence

Our univariate test is for the null hypothesis

$$X_{t,n,i,a} \perp\!\!\!\perp Y_{t,n,j,b} \mid \mathbf{Z}_{t,n} \text{ for all } t \in \mathcal{T}_n, \quad (1)$$

for a single dimension/time-offset tuple $(i, j, a, b) \in \mathcal{D}_n$. If domain knowledge suggests that we can restrict \mathcal{P}_n to consist of distributions in which the conditional dependencies are time-invariant, then we can use the alternative hypothesis

$$X_{t,n,i,a} \not\perp\!\!\!\perp Y_{t,n,j,b} \mid \mathbf{Z}_{t,n} \text{ for all } t \in \mathcal{T}_n. \quad (2)$$

We begin by focusing on time-invariant conditional independence relationships and address the time-varying case at the conclusion of this subsection.

Consider the following forecasting example in the univariate setting with $d_X = 1$, $d_Y = 1$, $\mathbf{d}_Z \geq 1$.

Example 2.1 (Univariate test for time series forecasting). *Suppose we are interested in determining whether the current value (time-offset $a = 0$) of a signal (dimension $i = 1$) is relevant for forecasting a target (dimension $j = 1$) seven time steps into the future (time-offset $b = 7$) after accounting for the existing forecasting signals $\mathbf{Z}_{t,n}$. Note that $\mathbf{Z}_{t,n}$ can consist of current values and lags of each of the forecasting signals. In this case, we would use the univariate version of our test with the null hypothesis*

$$X_{t,n,1,0} \perp\!\!\!\perp Y_{t,n,1,7} \mid \mathbf{Z}_{t,n} \text{ for all } t \in \mathcal{T}_n.$$

Naturally, our conditional independence test can be paired with multiple testing procedures so that we can conduct several of these univariate tests while controlling the false discovery rate. For example, we can individually test for conditional independence between $X_{t,n,i,a}$ and $Y_{t,n,j,b}$ given $\mathbf{Z}_{t,n}$ for any combination of forecasting horizons $b \in \{7, 14, 21, 28\}$, lags $a \in \{0, -7, -14\}$, targets $j \in \{1, 2, 3\}$, and signals $i \in \{1, 2, 3, 4, 5\}$.

The following null hypothesis of “no causal effect at all times” is pertinent to the growing literature on causal inference for time series [143, 127, 139, 140, 137], particularly for the setting in which just one realization of a nonstationary process is available.

Example 2.2 (Univariate test for time series causal inference). *Suppose we are interested in determining whether the current value (time-offset $a = 0$) of a continuous treatment (dimension $i = 1$) has any causal effect on an outcome of interest (dimension $j = 1$) one time step into the future (time-offset $b = 1$) after accounting for the confounders $\mathbf{Z}_{t,n}$. Crucially, we assume there are no unobserved confounders. In this case, we would use the univariate version of our test with the null hypothesis*

$$X_{t,n,1,0} \perp\!\!\!\perp Y_{t,n,1,1} \mid \mathbf{Z}_{t,n} \text{ for all } t \in \mathcal{T}_n.$$

If we do not require a p-value for each conditional independence relationship, then we can use the multivariate version of our test. In this case, we use the null hypothesis

$$X_{t,n,i,a} \perp\!\!\!\perp Y_{t,n,j,b} \mid \mathbf{Z}_{t,n} \text{ for all } t \in \mathcal{T}_n, \text{ for all } (i, j, a, b) \in \mathcal{D}_n. \quad (3)$$

By grouping together highly correlated dimensions and consecutive time-offsets, one can often construct a multivariate test that is more powerful than a univariate test based on a single dimension/time-offset tuple. Note that when using the multivariate test, different alternative hypotheses can be used depending on whether it is reasonable to restrict \mathcal{P}_n to consist of distributions in which the conditional dependencies are dimension-invariant. As we will explain next, there are many other situations in which grouping together related time series and using our multivariate test is useful.

Suppose we have time series data from nearby countries, cities, or sensors, as is common in the earth sciences, social sciences, and epidemiology. Write $X_{t,n,i,a}^\ell, Y_{t,n,j,b}^\ell, \mathbf{Z}_{t,n}^\ell$ to denote $X_{t,n,i,a}, Y_{t,n,j,b}, \mathbf{Z}_{t,n}$ at index $\ell \in \mathcal{L}_n$, where \mathcal{L}_n is an index set (e.g. different locations). We can often gain power by grouping together the time series in \mathcal{L}_n and using the multivariate test with the null hypothesis

$$X_{t,n,i,a}^\ell \perp\!\!\!\perp Y_{t,n,j,b}^\ell \mid \mathbf{Z}_{t,n}^\ell \text{ for all } t \in \mathcal{T}_n, \text{ for all } \ell \in \mathcal{L}_n, \quad (4)$$

for a single dimension/time-offset tuple $(i, j, a, b) \in \mathcal{D}_n$. Crucially, we allow each of the time series to be correlated with one another across \mathcal{L}_n and to have different distributions. Note that this is the same as the null hypothesis (1) but for all indices $\ell \in \mathcal{L}_n$. We can use the alternative hypothesis

$$X_{t,n,i,a}^\ell \not\perp\!\!\!\perp Y_{t,n,j,b}^\ell \mid \mathbf{Z}_{t,n}^\ell \text{ for all } t \in \mathcal{T}_n, \text{ for some } \ell \in \mathcal{L}_n, \quad (5)$$

since we have restricted the collection of distributions \mathcal{P}_n to consist of those in which the conditional dependencies are time-invariant. If we can further restrict \mathcal{P}_n so that it consists of distributions with time-invariant *and* index-invariant conditional dependencies, then we can use the alternative hypothesis

$$X_{t,n,i,a}^\ell \not\perp\!\!\!\perp Y_{t,n,j,b}^\ell \mid \mathbf{Z}_{t,n}^\ell \text{ for all } t \in \mathcal{T}_n, \text{ for all } \ell \in \mathcal{L}_n. \quad (6)$$

The following example is about forecasting a group of time series with $d_X = 1, d_Y = 1, \mathbf{d}_Z \geq 1$, and $|\mathcal{L}_n| > 1$. Similar hypotheses arise in causal inference with groups of time series (i.e. multivariate analogies of Example 2.2); see Section 6 of Wahl, Ninad, and Runge [162]. In many cases, the multivariate test used here will have more power than the univariate test used in Example 2.1.

Example 2.3 (Multivariate test for forecasting a group of time series). *As in Example 2.1, we are interested in forecasting a target seven time steps ahead. We want to determine whether the current value of a new forecasting signal is relevant or not after accounting for the existing forecasting signals. As before, we have dimensions $i = 1, j = 1$ and time-offsets $a = 0, b = 7$. However, now we have access to the same set of forecasting signals and targets at each location index $\ell \in \mathcal{L}_n$. In this case, we would use the multivariate version of our test with the null hypothesis*

$$X_{t,n,1,0}^\ell \perp\!\!\!\perp Y_{t,n,1,7}^\ell \mid \mathbf{Z}_{t,n}^\ell \text{ for all } t \in \mathcal{T}_n, \text{ for all } \ell \in \mathcal{L}_n.$$

Going forward, we suppress the superscript $\ell \in \mathcal{L}_n$ and revert back to the original notation (i.e. from $X_{t,n,i,a}^\ell, Y_{t,n,j,b}^\ell, \mathbf{Z}_{t,n}^\ell$ to $X_{t,n,i,a}, Y_{t,n,j,b}, \mathbf{Z}_{t,n}$). Note that this superscript can always be ignored outside of the “groups of time series” context, such as when there is only one index (e.g. one location).

To deal with the problem of time-varying conditional dependencies, we suggest modeling the conditional dependencies as though they are stable during certain time windows. If the breakpoints separating these time windows are known, then we can simply use our conditional independence test on each of these time windows and use multiple testing procedures to control the false discovery rate. However, this becomes more challenging if the breakpoints are unknown. In future work, we will develop a procedure for identifying time windows during which stable conditional dependencies hold while controlling the false discovery rate. That way, we can focus this manuscript on the main testing procedure. In Subsection A.5, we discuss how to test for time-varying conditional independence relationships at particular points in time by using the framework of locally stationary processes.

2.3 Time-varying regression functions

For a fixed sample size $n \in \mathbb{N}$, distribution $P \in \mathcal{P}_n$, time $t \in \mathcal{T}_n$ and dimension/time-offset tuple $(i, j, a, b) \in \mathcal{D}_n$, we can always decompose

$$X_{t,n,i,a} = f_{P,t,n,i,a}(\mathbf{Z}_{t,n}) + \varepsilon_{P,t,n,i,a}, \quad Y_{t,n,j,b} = g_{P,t,n,j,b}(\mathbf{Z}_{t,n}) + \xi_{P,t,n,j,b}, \quad (7)$$

where $f_{P,t,n,i,a}(\mathbf{z}) = \mathbb{E}_P(X_{t,n,i,a} | \mathbf{Z}_{t,n} = \mathbf{z})$ and $g_{P,t,n,j,b}(\mathbf{z}) = \mathbb{E}_P(Y_{t,n,j,b} | \mathbf{Z}_{t,n} = \mathbf{z})$ are the time-varying regression functions. The observed processes and error processes can all be nonstationary processes; see Section 3 for the details. Denote the product of errors at time t by

$$R_{P,t,n,m} = \varepsilon_{P,t,n,i,a} \xi_{P,t,n,j,b},$$

for $m = (i, j, a, b) \in \mathcal{D}_n$.

Next, let $\hat{f}_{t,n,i,a}$ and $\hat{g}_{t,n,j,b}$ be estimates of $f_{P,t,n,i,a}$ and $g_{P,t,n,j,b}$ created by time-varying nonlinear regressions of $(X_{t,n,i,a})_{t \in \mathcal{T}_n}$ on $(\mathbf{Z}_{t,n})_{t \in \mathcal{T}_n}$ and $(Y_{t,n,j,b})_{t \in \mathcal{T}_n}$ on $(\mathbf{Z}_{t,n})_{t \in \mathcal{T}_n}$, respectively. Let

$$\hat{\varepsilon}_{t,n,i,a} = X_{t,n,i,a} - \hat{f}_{t,n,i,a}(\mathbf{Z}_{t,n}), \quad \hat{\xi}_{t,n,j,b} = Y_{t,n,j,b} - \hat{g}_{t,n,j,b}(\mathbf{Z}_{t,n}),$$

be the corresponding residuals, and denote the product of these residuals at time t by

$$\hat{R}_{t,n,m} = \hat{\varepsilon}_{t,n,i,a} \hat{\xi}_{t,n,j,b}, \quad (8)$$

for $m = (i, j, a, b) \in \mathcal{D}_n$. Let $\hat{\mathbf{R}}_{t,n} = (\hat{R}_{t,n,m})_{m \in \mathcal{D}_n}$ be the high-dimensional vector process containing all the residual products for all dimension/time-offset combinations in \mathcal{D}_n .

2.4 Main ideas of our work and the algorithm

We show how to conduct tests for each of the null hypotheses from Subsection 2.2. Our proposed conditional independence tests are asymptotically valid as $n \rightarrow \infty$, *uniformly* over a large collection of distributions for which the null hypothesis holds. As mentioned in Section 1, our conditional independence test is inspired by the generalized covariance measure (GCM) test from Shah and Peters [149]. The original GCM test is based on expected conditional covariances between iid random variables, whereas our proposed dynamic GCM (dGCM) test is based on time-varying expected conditional covariances between nonstationary nonlinear processes. While our dGCM test shares similarities with the original GCM test from Shah and Peters [149], adapting it to this complex setting necessitates the use of advanced theoretical tools for nonstationary nonlinear processes, many of which have only been developed in recent years. Now, let us briefly summarize the main ideas behind the univariate version of the original GCM test.

For this paragraph, momentarily redefine X, Y to be two random variables and Z to be a random vector. The GCM test is based on the “weak” conditional independence criterion of Daudin [41], which states that if $X \perp\!\!\!\perp Y \mid Z$, then $\mathbb{E}_P[\phi(X, Z)\varphi(Y, Z)] = 0$ for all functions $\phi \in L^2_{X,Z}$ and $\varphi \in L^2_{Y,Z}$ such that $\mathbb{E}_P[\phi(X, Z) \mid Z] = 0$ and $\mathbb{E}_P[\varphi(Y, Z) \mid Z] = 0$. Thus, under the null hypothesis of conditional independence, the expectation of the products of errors $\mathbb{E}_P(\varepsilon\xi)$ from the regressions $X = \phi(Z) + \varepsilon$ and $Y = \varphi(Z) + \xi$, or equivalently the expected conditional covariance $\mathbb{E}_P[\text{Cov}_P(X, Y | Z)]$, is equal to

zero. As discussed in Shah and Peters [149], this can be seen as a generalization of the fact that the partial correlation coefficient, defined as the correlation between the residuals of linear regressions of X on Z and Y on Z , is equal to zero if and only if $X \perp\!\!\!\perp Y \mid Z$ when (X, Y, Z) are jointly Gaussian. The GCM test from Shah and Peters [149] is based on the normalized sum of the products of residuals from the regressions of X on Z and Y on Z .

Now, let us translate the “weak” conditional independence criterion of Daudin [41] into our setting using the notation from Subsection 2.1. If $X_{t,n,i,a} \perp\!\!\!\perp Y_{t,n,j,b} \mid \mathbf{Z}_{t,n}$, then

$$\mathbb{E}_P[\phi(X_{t,n,i,a}, \mathbf{Z}_{t,n})\varphi(Y_{t,n,j,b}, \mathbf{Z}_{t,n})] = 0,$$

for all functions $\phi \in L^2_{X_{t,n,i,a}, \mathbf{Z}_{t,n}}$ and $\varphi \in L^2_{Y_{t,n,j,b}, \mathbf{Z}_{t,n}}$ such that $\mathbb{E}_P[\phi(X_{t,n,i,a}, \mathbf{Z}_{t,n}) \mid \mathbf{Z}_{t,n}] = 0$ and $\mathbb{E}_P[\varphi(Y_{t,n,j,b}, \mathbf{Z}_{t,n}) \mid \mathbf{Z}_{t,n}] = 0$. Hence, under the null hypothesis, the expected conditional covariances,

$$\rho_{P,t,n,m} = \mathbb{E}_P[\text{Cov}_P(X_{t,n,i,a}, Y_{t,n,j,b} \mid \mathbf{Z}_{t,n})],$$

are always equal to zero for the dimension/time-offset combination $m = (i, j, a, b) \in \mathcal{D}_n$. Equivalently, the mean of the error products $\mathbb{E}_P(R_{P,t,n,m})$ from the time-varying nonlinear regressions of $(X_{t,n,i,a})_{t \in \mathcal{T}_n}$ on $(\mathbf{Z}_{t,n})_{t \in \mathcal{T}_n}$ and $(Y_{t,n,j,b})_{t \in \mathcal{T}_n}$ on $(\mathbf{Z}_{t,n})_{t \in \mathcal{T}_n}$ from Subsection 2.3 will always be zero under the null. This can be seen as a generalization of the partial correlation coefficient being equal to zero under conditional independence in the linear-Gaussian time series context; see the related discussion about Gaussian graphical models for nonstationary time series in Basu and Rao [11].

Crucially, the expected conditional covariances $\rho_{P,t,n,m}$ can be zero at all times, even under alternatives in which the corresponding conditional dependencies always hold. Consequently, we can only hope to have power against alternatives in which the time-varying expected conditional covariances $\rho_{P,t,n,m}$ are non-zero for at least *some* times. Hence, our test statistic (10) is designed to detect non-zero covariances between the errors $\varepsilon_{P,t,n,i,a}$ and $\xi_{P,t,n,j,b}$ at any point in time along the path.

We use a bootstrap-based testing procedure which appeals to the strong Gaussian approximation in Section B. The key ingredient of this bootstrap procedure is the time-varying covariance structure of the approximating nonstationary Gaussian process. Define the rolling window estimate of the time-varying covariance matrices of the vectors of error products by

$$\hat{\Sigma}_{t,n}^R = \frac{1}{L_n} \left(\sum_{s=t-L_n+1}^t \hat{R}_{s,n} \right)^{\otimes 2}, \quad (9)$$

where $L_n \in \mathbb{N}$ is a lag-window size parameter and the outer product is denoted by $v^{\otimes 2} = vv^T$. In the univariate case with dimensions $d_X = 1$, $d_Y = 1$, $d_Z \geq 1$ and time-offsets $A = \{a\}$, $B = \{b\}$, we use the rolling-window estimate of the time-varying variances of the error products

$$\hat{\sigma}_{t,n}^R = \frac{1}{L_n} \left(\sum_{s=t-L_n+1}^t \hat{R}_{s,n,m} \right)^2,$$

where $m = (1, 1, a, b) \in \mathcal{D}_n$. We postpone the details about these covariances until Subsection 3.5, and we discuss how to select L_n in Subsection 5.1.

Next, we introduce our test statistic, which is based on the process of residual products from the time-varying regressions of $(X_{t,n,i,a})_{t \in \mathcal{T}_n}$ on $(\mathbf{Z}_{t,n})_{t \in \mathcal{T}_n}$ and $(Y_{t,n,j,b})_{t \in \mathcal{T}_n}$ on $(\mathbf{Z}_{t,n})_{t \in \mathcal{T}_n}$. Define the set of times

$$\mathcal{T}_{n,L} = \{L_n + \mathbb{T}_n^- - 1, \dots, \mathbb{T}_n^+ - 1, \mathbb{T}_n^+\},$$

and denote its cardinality by $T_{n,L} = |\mathcal{T}_{n,L}|$. Denote the entire process containing the residual products by $\hat{\mathbf{R}}_n = (\hat{\mathbf{R}}_{t,n})_{t \in \mathcal{T}_{n,L}}$. The test statistic based on the maximum ℓ_p -norm ($p \geq 2$) achieved by the partial sum process is given by

$$S_{n,p}(\hat{\mathbf{R}}_n) = \max_{s \in \mathcal{T}_{n,L}} \left\| \frac{1}{\sqrt{T_{n,L}}} \sum_{t \leq s} \hat{\mathbf{R}}_{t,n} \right\|_p. \quad (10)$$

Using this test statistic, we will reject the null hypothesis of conditional independence if the ℓ_p norm of the partial sum process of residual products ever becomes “too large” *at any point in time along*

the path. For example, we can use the ℓ_∞ -type or ℓ_2 -type test statistics

$$S_{n,\infty}(\hat{\mathbf{R}}_n) = \max_{s \in \mathcal{T}_{n,L}} \left\| \frac{1}{\sqrt{T_{n,L}}} \sum_{t \leq s} \hat{\mathbf{R}}_{t,n} \right\|_\infty, \quad S_{n,2}(\hat{\mathbf{R}}_n) = \max_{s \in \mathcal{T}_{n,L}} \left\| \frac{1}{\sqrt{T_{n,L}}} \sum_{t \leq s} \hat{\mathbf{R}}_{t,n} \right\|_2,$$

to achieve high power against either sparse or dense alternatives, respectively. In the univariate case with dimensions $d_X = 1$, $d_Y = 1$, $d_Z \geq 1$ and time-offsets $A = \{a\}$, $B = \{b\}$, the test statistic reduces to the absolute value of the partial sum process of residual products

$$S_n(\hat{R}_{n,m}) = \max_{s \in \mathcal{T}_{n,L}} \left| \frac{1}{\sqrt{T_{n,L}}} \sum_{t \leq s} \hat{R}_{t,n,m} \right|,$$

where $m = (1, 1, a, b) \in \mathcal{D}_n$ and $\hat{R}_{n,m} = (\hat{R}_{t,n,m})_{t \in \mathcal{T}_{n,L}}$. See Subsections A.1 and A.5 for discussions of alternative test statistics, namely those based on the full sum and those employing kernel smoothing.

The dGCM test is given by Algorithm 1 below. The main steps are time-varying nonlinear regression, time-varying covariance estimation, and a bootstrap procedure justified by the distribution-uniform strong Gaussian approximation in Section B. We present the testing procedure in the multivariate setting, and the algorithm for the univariate setting is obtained by replacing $S_{n,p}(\cdot)$, $\hat{\Sigma}_{t,n}^R$, $\hat{\mathbf{R}}_{t,n}$, $\check{\mathbf{R}}_n$, $\check{\mathbf{R}}_{t,n}^{(r)}$, $\check{\mathbf{R}}_n^{(r)}$ with $S_n(\cdot)$, $\hat{\sigma}_{t,n}^R$, $\hat{R}_{t,n}$, \check{R}_n , $\check{R}_{t,n}^{(r)}$, $\check{R}_n^{(r)}$, respectively.

Algorithm 1 The dynamic generalized covariance measure (dGCM) test

- 1: **Input:** Dimensions and time-offsets of interest \mathcal{D}_n , time points \mathcal{T}_n , data $(X_{t,n}, Y_{t,n}, Z_{t,n})_{t \in \mathcal{T}_n}$, test statistic $S_{n,p}(\cdot)$, α for the significance level, α' for the quantile $\hat{q}_{1-\alpha'}^{\text{boot}}$, number of simulations s
 - 2: **for** each time $t \in \mathcal{T}_n$ and dimension/time-offset tuple $m = (i, j, a, b) \in \mathcal{D}_n$ **do**
 - 3: Obtain estimates $\hat{f}_{t,n,i,a}$ and $\hat{g}_{t,n,j,b}$ of the time-varying regression functions from (7)
 - 4: Calculate the product of residuals $\hat{R}_{t,n,m} = \hat{e}_{t,n,i,a} \hat{\xi}_{t,n,j,b}$ from (8)
 - 5: **end for**
 - 6: Select the lag-window size L_n for covariance estimation according to Subsection 5.1
 - 7: **for** each time $t \in \mathcal{T}_{n,L}$ **do**
 - 8: Calculate the rolling-window estimates $\hat{\Sigma}_{t,n}^R$ of the time-varying covariance matrices from (9)
 - 9: **end for**
 - 10: **for** each simulation $r = 1, \dots, s$ **do**
 - 11: **for** each time $t \in \mathcal{T}_{n,L}$ **do**
 - 12: Simulate independent Gaussian random vectors $\check{\mathbf{R}}_{t,n}^{(r)} \sim \mathcal{N}(0, \hat{\Sigma}_{t,n}^R)$
 - 13: **end for**
 - 14: Calculate the test statistic $S_{n,p}(\check{\mathbf{R}}_n^{(r)})$ from (10) using the Gaussian process $\check{\mathbf{R}}_n^{(r)} = (\check{\mathbf{R}}_{t,n}^{(r)})_{t \in \mathcal{T}_{n,L}}$
 - 15: **end for**
 - 16: Calculate the $1 - \alpha'$ empirical quantile $\hat{q}_{1-\alpha'}^{\text{boot}}$ of $(S_{n,p}(\check{\mathbf{R}}_n^{(r)}))_{r=1}^s$
 - 17: Calculate the test statistic $S_{n,p}(\hat{\mathbf{R}}_n)$ from (10) using the residual products $\hat{\mathbf{R}}_n = (\hat{\mathbf{R}}_{t,n})_{t \in \mathcal{T}_{n,L}}$
 - 18: **if** $S_{n,p}(\hat{\mathbf{R}}_n) > \hat{q}_{1-\alpha'}^{\text{boot}}$ **then**
 - 19: Reject the null hypothesis at the level α
 - 20: **end if**
 - 21: **Output:** Decision to either reject or fail to reject the null hypothesis at the level α
-

We investigate the finite sample performance of the dGCM test in Section 5. While the dGCM test can be used with any black-box estimator of the time-varying regression functions from (7), we use the sieve estimator from Section 4 in our simulations. As with the GCM test from Shah and Peters [149], the Type-I error control guarantee for the dGCM test is asymptotic; see Theorem 3.1 for the details. Our simulations indicate that the dGCM test maintains the level if the sample size is large enough to reliably estimate the time-varying regression functions. With large sample sizes, we can straightforwardly reject the null hypothesis when the test statistic exceeds the quantile $\hat{q}_{0.95}^{\text{boot}}$ from Algorithm 1 to obtain an approximately level 0.05 test. On the other hand, with small sample sizes, we may be conservative and only reject the null hypothesis at level 0.05 if the test statistic exceeds the quantile $\hat{q}_{0.975}^{\text{boot}}$ from Algorithm 1. For more discussion about the simulation results, see Subsection 5.3.

3 Assumptions and the Theoretical Result for dGCM

In this section, we introduce a triangular array framework for high-dimensional nonstationary nonlinear processes. Specifically, our framework is designed to enable hypothesis testing based on residuals formed from time-varying nonlinear regression estimates. It provides a foundational template for future research in causal inference, graphical modeling, and variable selection for this setting.

We allow the processes to have long-range temporal dependence and very complicated forms of nonstationarity which can be both abrupt and smooth. We control the temporal dependence and nonstationarity of the processes *uniformly* over collections of distributions by employing versions of the functional dependence measure of Wu [170] and the total variation-type nonstationarity condition of Mies and Steland [109]; see Assumptions 3.5 and 3.6 for the details. These distribution-uniform assumptions are needed for the uniform level guarantee in Theorem 3.1, which is our main theoretical result. We discuss the importance of this uniform guarantee in Subsection 5.3.

The framework we introduce in this section nests several well-studied classes of processes. In Section 4, we show how our framework nests a class of nonstationary processes called *locally stationary processes*, which allows for smooth changes over time. We refer interested readers to Dahlhaus [37] and Dahlhaus, Richter, and Wu [40] for more information about the linear and nonlinear cases, respectively. In Subsection A.6, we explain how our framework also nests a more general class of nonstationary processes called *piecewise locally stationary processes*. This class extends the framework for locally stationary processes by permitting both smooth changes and abrupt breakpoints. Naturally, our framework includes the class of stationary processes as a special case when there is no nonstationarity. In Subsection A.4, we discuss how our test can be simplified when stationarity can be assumed. Similarly, the fundamental setting of iid sequences arises as another special case when there is neither nonstationarity nor temporal dependence.

Lastly, we discuss how our general triangular array framework is compatible with even more types of nonstationary processes in the Appendix. In Subsection A.2, we consider a class of nonstationary processes called *cyclostationary processes* which exhibit repetition over time. Also, in Subsection A.3, we explain how our framework can leverage black-box simulators for the multivariate process (X, Z) by using simulation-and-regression techniques.

3.1 Nonstationary observed processes

In this subsection, we introduce the so-called “causal representations” of the processes. Specifically, we view each dimension of the observed sequence $(X_{t,n}, Y_{t,n}, Z_{t,n})_{t \in [n]}$ as the outputs of a time-varying nonlinear function that is given a sequence of iid inputs. This type of representation has a long history, tracing back to at least Rosenblatt [135] and Wiener [167], though its importance for the statistical analysis of time series was first elucidated by Wu [170]. What follows is most similar to the framework for high-dimensional nonstationary nonlinear processes from Mies and Steland [109], which in turn builds on the framework from Zhou and Wu [190]. For the following assumption, let

$$\mathcal{H}_t^X = (\eta_t^X, \eta_{t-1}^X, \dots), \mathcal{H}_t^Y = (\eta_t^Y, \eta_{t-1}^Y, \dots), \mathcal{H}_t^Z = (\eta_t^Z, \eta_{t-1}^Z, \dots),$$

where $(\eta_t^X, \eta_t^Y, \eta_t^Z)_{t \in \mathbb{Z}}$ is a sequence of iid random vectors. Denote the dimensions of $\eta_t^X = \eta_{t,n}^X$, $\eta_t^Y = \eta_{t,n}^Y$, $\eta_t^Z = \eta_{t,n}^Z$ respectively by $d_X^\eta = d_{X,n}^\eta$, $d_Y^\eta = d_{Y,n}^\eta$, $d_Z^\eta = d_{Z,n}^\eta$, which can change with n .

Assumption 3.1 (Causal representations of the observed processes). *Assume that for each time $t \in \mathcal{T}_n$ we can represent each dimension of each of the observed processes as the output of an evolving nonlinear system that was given a sequence of iid inputs:*

$$X_{t,n,i} = G_{t,n,i}^X(\mathcal{H}_t^X), Y_{t,n,j} = G_{t,n,j}^Y(\mathcal{H}_t^Y), Z_{t,n,k} = G_{t,n,k}^Z(\mathcal{H}_t^Z).$$

For each $n \in \mathbb{N}$, $(i, j, a, b) \in \mathcal{D}_n$, $t \in \mathcal{T}_n$, we assume that $G_{t,n,i}^X(\cdot)$, $G_{t,n,j}^Y(\cdot)$, $G_{t,n,k}^Z(\cdot)$ are each measurable functions from $(\mathbb{R}^{d_X^\eta})^\infty$, $(\mathbb{R}^{d_Y^\eta})^\infty$, $(\mathbb{R}^{d_Z^\eta})^\infty$ to \mathbb{R} — where we endow $(\mathbb{R}^{d_X^\eta})^\infty$, $(\mathbb{R}^{d_Y^\eta})^\infty$, $(\mathbb{R}^{d_Z^\eta})^\infty$ with the σ -algebra generated by all finite projections — such that $G_{t,n,i}^X(\mathcal{H}_s^X)$, $G_{t,n,j}^Y(\mathcal{H}_s^Y)$, $G_{t,n,k}^Z(\mathcal{H}_s^Z)$ are each well-defined random variables for each $s \in \mathbb{Z}$ and $(G_{t,n,i}^X(\mathcal{H}_s^X))_{s \in \mathbb{Z}}$, $(G_{t,n,j}^Y(\mathcal{H}_s^Y))_{s \in \mathbb{Z}}$, $(G_{t,n,k}^Z(\mathcal{H}_s^Z))_{s \in \mathbb{Z}}$ are each stationary ergodic processes.

To simplify the notation, we have not defined the input sequences for the observed processes separately for each dimension. Without loss of generality, we can define the measurable functions $G_{t,n,i}^X(\cdot)$, $G_{t,n,j}^Y(\cdot)$, $G_{t,n,k}^Z(\cdot)$ and the inputs η_t^X , η_t^Y , η_t^Z so that each dimension of the observed processes can have idiosyncratic inputs.

We will introduce several more causal representations throughout this paper. Let us state some properties that all causal representations will have to avoid repeating the same ideas each time. The causal representations will all be measurable functions on $(\mathbb{R}^{d^n})^\infty$ for some $d^n \in \mathbb{N}$, where we will always endow $(\mathbb{R}^{d^n})^\infty$ with the σ -algebra generated by all finite projections. The causal mechanism at a particular time $t \in \mathcal{T}_n$ with the input sequence up to a particular $s \in \mathbb{Z}$ is a well-defined random variable or vector. Similarly, the process induced by considering the input sequence up to each $s \in \mathbb{Z}$ with a fixed causal mechanism is a stationary ergodic process, as in Assumption 3.1.

In view of Assumption 3.1, we have the following causal representations for the observed processes with all dimensions

$$\begin{aligned} X_{t,n} &= G_{t,n}^X(\mathcal{H}_t^X) = (G_{t,n,i}^X(\mathcal{H}_t^X))_{i \in [d_X]}, \\ Y_{t,n} &= G_{t,n}^Y(\mathcal{H}_t^Y) = (G_{t,n,j}^Y(\mathcal{H}_t^Y))_{j \in [d_Y]}, \\ Z_{t,n} &= G_{t,n}^Z(\mathcal{H}_t^Z) = (G_{t,n,k}^Z(\mathcal{H}_t^Z))_{k \in [d_Z]}. \end{aligned}$$

Also, we have causal representations for each of the dimensions $i \in [d_X]$, $j \in [d_Y]$, $k \in [d_Z]$ with time-offsets $a \in A_i$, $b \in B_j$, $c \in C_k$

$$\begin{aligned} X_{t,n,i,a} &= G_{t,n,i,a}^X(\mathcal{H}_{t,a}^X) = G_{t+a,n,i}^X(\mathcal{H}_{t+a}^X), \\ Y_{t,n,j,b} &= G_{t,n,j,b}^Y(\mathcal{H}_{t,b}^Y) = G_{t+b,n,j}^Y(\mathcal{H}_{t+b}^Y), \\ Z_{t,n,k,c} &= G_{t,n,k,c}^Z(\mathcal{H}_{t,c}^Z) = G_{t+c,n,k}^Z(\mathcal{H}_{t+c}^Z), \end{aligned}$$

where $\mathcal{H}_{t,a}^X = (\eta_{t+a}^X, \eta_{t-1+a}^X, \dots)$, $\mathcal{H}_{t,b}^Y = (\eta_{t+b}^Y, \eta_{t-1+b}^Y, \dots)$, and $\mathcal{H}_{t,c}^Z = (\eta_{t+c}^Z, \eta_{t-1+c}^Z, \dots)$. We can then write the causal representation of the vectors with all dimensions and time-offsets as

$$\begin{aligned} \mathbf{X}_{t,n} &= \mathbf{G}_{t,n}^X(\mathcal{H}_t^X) = (G_{t,n,i,a}^X(\mathcal{H}_{t,a}^X))_{i \in [d_X], a \in A_i}, \\ \mathbf{Y}_{t,n} &= \mathbf{G}_{t,n}^Y(\mathcal{H}_t^Y) = (G_{t,n,j,b}^Y(\mathcal{H}_{t,b}^Y))_{j \in [d_Y], b \in B_j}, \\ \mathbf{Z}_{t,n} &= \mathbf{G}_{t,n}^Z(\mathcal{H}_t^Z) = (G_{t,n,k,c}^Z(\mathcal{H}_{t,c}^Z))_{k \in [d_Z], c \in C_k}, \end{aligned}$$

where $\mathcal{H}_t^X = (\eta_t^X, \eta_{t-1}^X, \dots)$, $\mathcal{H}_t^Y = (\eta_t^Y, \eta_{t-1}^Y, \dots)$, $\mathcal{H}_t^Z = (\eta_t^Z, \eta_{t-1}^Z, \dots)$, and $\eta_t^X = \eta_{t+a_{\max}}^X$, $\eta_t^Y = \eta_{t+b_{\max}}^Y$, $\eta_t^Z = \eta_{t+c_{\max}}^Z$.

Let Ω be a sample space, \mathcal{B} the Borel sigma-algebra, and (Ω, \mathcal{B}) a measurable space. For fixed $n \in \mathbb{N}$, let (Ω, \mathcal{B}) be equipped with a family of probability measures $(\mathbb{P}_P)_{P \in \mathcal{P}_n}$ so that the joint distribution of the nonlinear stochastic systems

$$(G_{t,n}^X(\mathcal{H}_s^X))_{t \in [n], s \in \mathbb{Z}}, (G_{t,n}^Y(\mathcal{H}_s^Y))_{t \in [n], s \in \mathbb{Z}}, (G_{t,n}^Z(\mathcal{H}_s^Z))_{t \in [n], s \in \mathbb{Z}}$$

under \mathbb{P}_P is $P \in \mathcal{P}_n$, where the collection of distributions \mathcal{P}_n can change with n . The family of probability measures $(\mathbb{P}_P)_{P \in \mathcal{P}_n}$ is defined with respect to the same measurable space (Ω, \mathcal{B}) , but need not have the same dominating measure. Denote the family of probability spaces by $(\Omega, \mathcal{B}, \mathbb{P}_P)_{P \in \mathcal{P}_n}$ and a sequence of such families of probability spaces by $((\Omega, \mathcal{B}, \mathbb{P}_P)_{P \in \mathcal{P}_n})_{n \in \mathbb{N}}$.

For a given sample size $n \in \mathbb{N}$ and distribution $P \in \mathcal{P}_n$, let $\mathbb{E}_P(\cdot)$ denote the expectation of a random variable with distribution determined by P . Let $\mathbb{P}_P(E)$ denote the probability of an event $E \in \mathcal{B}$. We use the notation $o_{\mathcal{P}}(\cdot)$ and $O_{\mathcal{P}}(\cdot)$ in the same way that Shah and Peters [149] do, so we replicate their notation here. Let $(V_{P,n})_{n \in \mathbb{N}, P \in \mathcal{P}_n}$ be a family of sequences of random variables with distributions determined by $P \in \mathcal{P}_n$ for some collection of distributions \mathcal{P}_n which will be made clear from the context. We write $V_{P,n} = o_{\mathcal{P}}(1)$ to mean that for all $\epsilon > 0$, we have

$$\sup_{P \in \mathcal{P}_n} \mathbb{P}_P(|V_{P,n}| > \epsilon) \rightarrow 0.$$

Also, by $V_{P,n} = O_{\mathcal{P}}(1)$ we mean for all $\epsilon > 0$, there exists a constant $K > 0$ such that

$$\sup_{n \in \mathbb{N}} \sup_{P \in \mathcal{P}_n} \mathbb{P}_P(|V_{P,n}| > K) < \epsilon.$$

Let $(W_{P,n})_{n \in \mathbb{N}, P \in \mathcal{P}_n}$ be another family of sequences of random variables. By $V_{P,n} = o_{\mathcal{P}}(W_{P,n})$ we mean $V_{P,n} = W_{P,n}R_{P,n}$ and $R_{P,n} = o_{\mathcal{P}}(1)$, and by $V_{P,n} = O_{\mathcal{P}}(W_{P,n})$ we mean $V_{P,n} = W_{P,n}R_{P,n}$ and $R_{P,n} = O_{\mathcal{P}}(1)$.

In the next few subsections, we will state distribution-uniform assumptions with respect to a *generic* sequence of collections of distributions $(\mathcal{P}_n)_{n \in \mathbb{N}}$ for the observed processes. Let $\mathcal{P}_{0,n}^{\text{CI}}$ be a collection of distributions for the observed processes such that the null hypothesis is true, and let $(\mathcal{P}_{0,n}^{\text{CI}})_{n \in \mathbb{N}}$ be a sequence of such collections of distributions. In our main result, which we state as Theorem 3.1, we will assume that these distribution-uniform assumptions hold for a sequence of collections of distributions $(\mathcal{P}_{0,n}^*)_{n \in \mathbb{N}}$, where $\mathcal{P}_{0,n}^* \subset \mathcal{P}_{0,n}^{\text{CI}}$ for each $n \in \mathbb{N}$. That is, we will make these assumptions for a sequence of subcollections of distributions for which the global null hypothesis of conditional independence holds.

3.2 Prediction processes

We introduce causal representations for the prediction processes in this subsection. For each $n \in \mathbb{N}$, $t \in \mathcal{T}_n$, $(i, j, a, b) \in \mathcal{D}_n$, let $\eta_{t,n,i,a}^{\text{algo}}, \eta_{t,n,j,b}^{\text{algo}}$ be random variables that encode the (possible) stochasticity of the statistical learning algorithms. If the learning algorithms are not stochastic, then these random variables can be ignored without loss of generality. Going forward, we will suppress the dependence of the predictors on $\eta_{t,n,i,a}^{\text{algo}}, \eta_{t,n,j,b}^{\text{algo}}$ to simplify the notation.

Let $\mathfrak{D}_{t,n,i,a}^{\hat{f}}, \mathfrak{D}_{t,n,j,b}^{\hat{g}}$ be the datasets containing the observations used to form the predictors $\hat{f}_{t,n,i,a}, \hat{g}_{t,n,j,b}$, and let $\mathcal{H}_{t,a}^{\mathfrak{D}^{\hat{f}}}, \mathcal{H}_{t,b}^{\mathfrak{D}^{\hat{g}}}$ be the corresponding input sequences. For example, if only the observations in \mathcal{T}_n up to time $t \in \mathcal{T}_n$ are used to form the predictor $\hat{g}_{t,n,j,b}$, then $\mathfrak{D}_{t,n,j,b}^{\hat{g}} = (Y_{s,n,j,b}, \mathbf{Z}_{s,n})_{s \leq t}$ and $\mathcal{H}_{t,b}^{\mathfrak{D}^{\hat{g}}} = (\mathcal{H}_{t,b}^Y, \mathcal{H}_t^Z)$. Similarly, if all of the observations in \mathcal{T}_n are used (i.e. to time \mathbb{T}_n^+) to form the predictor $\hat{g}_{t,n,j,b}$, then $\mathfrak{D}_{t,n,j,b}^{\hat{g}} = (Y_{t,n,j,b}, \mathbf{Z}_{t,n})_{t \in \mathcal{T}_n}$ and $\mathcal{H}_{t,b}^{\mathfrak{D}^{\hat{g}}} = (\mathcal{H}_{\mathbb{T}_n^+}^Y, \mathcal{H}_{\mathbb{T}_n^+}^Z)$.

Denote the sets of times corresponding to $\mathfrak{D}_{t,n,i,a}^{\hat{f}}, \mathfrak{D}_{t,n,j,b}^{\hat{g}}$ by $\mathcal{T}_{t,n,i,a}^{\hat{f}}, \mathcal{T}_{t,n,j,b}^{\hat{g}}$, respectively, and let $T_{t,n,i,a}^{\hat{f}} = |\mathcal{T}_{t,n,i,a}^{\hat{f}}|$, $T_{t,n,j,b}^{\hat{g}} = |\mathcal{T}_{t,n,j,b}^{\hat{g}}|$ be the cardinalities. For each $n \in \mathbb{N}$, $t \in \mathcal{T}_n$ let $\mathcal{M}(\mathcal{Z}, \mathcal{Y}) \subseteq \mathcal{Y}^{\mathcal{Z}}$ and $\mathcal{M}(\mathcal{Z}, \mathcal{X}) \subseteq \mathcal{X}^{\mathcal{Z}}$, where $\mathcal{X} = \mathbb{R}$, $\mathcal{Y} = \mathbb{R}$, and $\mathcal{Z} = \mathbb{R}^{\mathbf{d}_Z}$. Note that \mathbf{d}_Z can grow with n as discussed in Subsection 2.1, although we suppress this in the notation.

Assumption 3.2 (Causal representations of the predictors). *For each $n \in \mathbb{N}$, $(i, j, a, b) \in \mathcal{D}_n$, assume that the sequences of statistical learning algorithms $\mathcal{A}_{n,i,a}^{\hat{f}} = (\mathcal{A}_{t,n,i,a}^{\hat{f}})_{t \in \mathcal{T}_n}$, $\mathcal{A}_{n,j,b}^{\hat{g}} = (\mathcal{A}_{t,n,j,b}^{\hat{g}})_{t \in \mathcal{T}_n}$ consist of the Borel measurable functions*

$$\mathcal{A}_{t,n,i,a}^{\hat{f}} : \begin{cases} (\mathcal{Z} \times \mathcal{X})^{T_{t,n,i,a}^{\hat{f}}} \rightarrow \mathcal{M}(\mathcal{Z}, \mathcal{X}) \\ \mathfrak{D}_{t,n,i,a}^{\hat{f}} \rightarrow \hat{f}_{t,n,i,a}, \end{cases}$$

and

$$\mathcal{A}_{t,n,j,b}^{\hat{g}} : \begin{cases} (\mathcal{Z} \times \mathcal{Y})^{T_{t,n,j,b}^{\hat{g}}} \rightarrow \mathcal{M}(\mathcal{Z}, \mathcal{Y}) \\ \mathfrak{D}_{t,n,j,b}^{\hat{g}} \mapsto \hat{g}_{t,n,j,b}, \end{cases}$$

such that the predictors have the causal representations

$$\begin{aligned} \hat{f}_{t,n,i,a} &= G_{t,n,i,a}^{\mathcal{A}^{\hat{f}}}(\mathcal{H}_{t,a}^{\mathfrak{D}^{\hat{f}}}), \\ \hat{g}_{t,n,j,b} &= G_{t,n,j,b}^{\mathcal{A}^{\hat{g}}}(\mathcal{H}_{t,b}^{\mathfrak{D}^{\hat{g}}}), \end{aligned}$$

in view of Assumption 3.1. $G_{t,n,i,a}^{\mathcal{A}^{\hat{f}}}(\cdot)$, $G_{t,n,j,b}^{\mathcal{A}^{\hat{g}}}(\cdot)$ are measurable functions so that $G_{t,n,i,a}^{\mathcal{A}^{\hat{f}}}(\mathcal{H}_{t,a}^{\mathfrak{D}^{\hat{f}}})$, $G_{t,n,j,b}^{\mathcal{A}^{\hat{g}}}(\mathcal{H}_{t,b}^{\mathfrak{D}^{\hat{g}}})$ are well-defined function-valued random variables.

We make the following assumption for the predictions and prediction errors for some sequence of collections of distributions $(\mathcal{P}_n)_{n \in \mathbb{N}}$.

Assumption 3.3 (Causal representations of the predictions and prediction errors). *Assume that the predictors $\hat{f}_{t,n,i,a}, \hat{g}_{t,n,j,b}$ are Borel measurable functions from $\mathbb{R}^{\mathbf{d}_Z}$ to \mathbb{R} such that for each $n \in \mathbb{N}$,*

$t \in \mathcal{T}_n$, $(i, j, a, b) \in \mathcal{D}_n$ we can represent the predictions and prediction errors as

$$\begin{aligned}\hat{f}_{t,n,i,a}(\mathbf{Z}_{t,n}) &= G_{t,n,i,a}^{\hat{f}}(\mathcal{H}_{t,a}^{\hat{f}}) = [\mathcal{A}_{t,n,i,a}^{\hat{f}}(X_{n,i,a}, \mathbf{Z}_n)](\mathbf{Z}_{t,n}), \\ \hat{g}_{t,n,j,b}(\mathbf{Z}_{t,n}) &= G_{t,n,j,b}^{\hat{g}}(\mathcal{H}_{t,b}^{\hat{g}}) = [\mathcal{A}_{t,n,j,b}^{\hat{g}}(Y_{n,j,b}, \mathbf{Z}_n)](\mathbf{Z}_{t,n}),\end{aligned}$$

and

$$\begin{aligned}\hat{w}_{P,t,n,i,a}^f &= G_{P,t,n,i,a}^{\hat{w}^f}(\mathcal{H}_{t,a}^{\hat{f}}) = f_{P,t,n,i,a}(\mathbf{Z}_{t,n}) - \hat{f}_{t,n,i,a}(\mathbf{Z}_{t,n}), \\ \hat{w}_{P,t,n,j,b}^g &= G_{P,t,n,j,b}^{\hat{w}^g}(\mathcal{H}_{t,b}^{\hat{g}}) = g_{P,t,n,j,b}(\mathbf{Z}_{t,n}) - \hat{g}_{t,n,j,b}(\mathbf{Z}_{t,n}),\end{aligned}$$

in view of Assumptions 3.1, 3.2, where the input sequences are

$$\mathcal{H}_{t,a}^{\hat{f}} = (\mathcal{H}_{t,a}^{\mathfrak{D}^{\hat{f}}}, \mathcal{H}_t^{\mathbf{Z}}), \quad \mathcal{H}_{t,b}^{\hat{g}} = (\mathcal{H}_{t,b}^{\mathfrak{D}^{\hat{g}}}, \mathcal{H}_t^{\mathbf{Z}}).$$

Also, assume that for all $n \in \mathbb{N}$, $t \in \mathcal{T}_n$, $(i, j, a, b) \in \mathcal{D}_n$ there exists some $q \geq 2$ such that

$$\sup_{P \in \mathcal{P}_n} \mathbb{E}_P(|\hat{w}_{P,t,n,i,a}^f|^q) < \infty, \quad \sup_{P \in \mathcal{P}_n} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g|^q) < \infty.$$

$G_{t,n,i,a}^{\hat{f}}(\cdot)$, $G_{P,t,n,i,a}^{\hat{w}^f}(\cdot)$ and $G_{t,n,j,b}^{\hat{g}}(\cdot)$, $G_{P,t,n,j,b}^{\hat{w}^g}(\cdot)$ are measurable functions such that $G_{t,n,i,a}^{\hat{f}}(\mathcal{H}_{t,a}^{\hat{f}})$, $G_{t,n,j,b}^{\hat{g}}(\mathcal{H}_{t,b}^{\hat{g}})$ and $G_{P,t,n,i,a}^{\hat{w}^f}(\mathcal{H}_{t,a}^{\hat{f}})$, $G_{P,t,n,j,b}^{\hat{w}^g}(\mathcal{H}_{t,b}^{\hat{g}})$ are well-defined real-valued random variables.

In view of Assumption 3.3, we have the following causal representation for all dimensions and time-offsets of the prediction errors

$$\begin{aligned}\hat{w}_{P,t,n}^f &= \mathbf{G}_{P,t,n}^{\hat{w}^f}(\mathcal{H}_t^{\hat{f}}) = (\hat{w}_{P,t,n,i,a}^f)_{i \in [d_X], a \in A_i}, \\ \hat{w}_{P,t,n}^g &= \mathbf{G}_{P,t,n}^{\hat{w}^g}(\mathcal{H}_t^{\hat{g}}) = (\hat{w}_{P,t,n,j,b}^g)_{j \in [d_Y], b \in B_j},\end{aligned}$$

where $\mathcal{H}_t^{\hat{f}} = (\mathcal{H}_{t,a}^{\hat{f}})_{a \in A}$ and $\mathcal{H}_t^{\hat{g}} = (\mathcal{H}_{t,b}^{\hat{g}})_{b \in B}$.

3.3 Nonstationary error processes

In this subsection, we will introduce the causal representations of the error processes of the high-dimensional nonstationary processes. For the next assumption, for each $a \in A$, $b \in B$, define the input sequences

$$\mathcal{H}_{t,a}^{\varepsilon} = (\eta_{t,a}^{\varepsilon}, \eta_{t,a-1}^{\varepsilon}, \dots), \quad \mathcal{H}_{t,b}^{\xi} = (\eta_{t,b}^{\xi}, \eta_{t,b-1}^{\xi}, \dots), \quad (11)$$

where $(\eta_{t,a}^{\varepsilon}, \eta_{t,b}^{\xi})_{t \in \mathbb{Z}}$ is a sequence of iid random vectors. For the following assumption, denote the dimension of $\eta_{t,a}^{\varepsilon} = \eta_{t,a,n}^{\varepsilon}$ by $d_{\varepsilon}^{\eta} = d_{\varepsilon,n}^{\eta}$, and the dimension of $\eta_{t,b}^{\xi} = \eta_{t,b,n}^{\xi}$ by $d_{\xi}^{\eta} = d_{\xi,n}^{\eta}$, both of which can change with n .

Assumption 3.4 (Causal representations of the error processes). *Assume that for each $n \in \mathbb{N}$, $P \in \mathcal{P}_n$, $(i, j, a, b) \in \mathcal{D}_n$, $t \in \mathcal{T}_n$, we can represent the error processes from Subsection 2.3 as*

$$\varepsilon_{P,t,n,i,a} = G_{P,t,n,i,a}^{\varepsilon}(\mathcal{H}_{t,a}^{\varepsilon}), \quad \xi_{P,t,n,j,b} = G_{P,t,n,j,b}^{\xi}(\mathcal{H}_{t,b}^{\xi}),$$

with $\mathbb{E}_P(\varepsilon_{P,t,n,i,a} | \mathcal{H}_t^{\hat{g}}) = 0$ and $\mathbb{E}_P(\xi_{P,t,n,j,b} | \mathcal{H}_t^{\hat{f}}) = 0$, where the input sequences $\mathcal{H}_t^{\hat{g}}$, $\mathcal{H}_t^{\hat{f}}$ are defined following Assumption 3.3. $G_{P,t,n,i,a}^{\varepsilon}(\cdot)$ and $G_{P,t,n,j,b}^{\xi}(\cdot)$ are measurable functions from $(\mathbb{R}^{d_{\varepsilon}^{\eta}})^{\infty}$ and $(\mathbb{R}^{d_{\xi}^{\eta}})^{\infty}$, respectively, to \mathbb{R} — where we endow $(\mathbb{R}^{d_{\varepsilon}^{\eta}})^{\infty}$ and $(\mathbb{R}^{d_{\xi}^{\eta}})^{\infty}$ with the σ -algebra generated by all finite projections — so that $G_{P,t,n,i,a}^{\varepsilon}(\mathcal{H}_{s,a}^{\varepsilon})$, $G_{P,t,n,j,b}^{\xi}(\mathcal{H}_{s,b}^{\xi})$ are well-defined random variables for each $s \in \mathbb{Z}$ and $(G_{P,t,n,i,a}^{\varepsilon}(\mathcal{H}_{s,a}^{\varepsilon}))_{s \in \mathbb{Z}}$, $(G_{P,t,n,j,b}^{\xi}(\mathcal{H}_{s,b}^{\xi}))_{s \in \mathbb{Z}}$ are stationary ergodic processes.

We have not defined the input sequences for the error processes separately for each dimension. Without loss of generality, the measurable functions $G_{P,t,n,i,a}^{\varepsilon}(\cdot)$, $G_{P,t,n,j,b}^{\xi}(\cdot)$ and inputs $\eta_{t,a}^{\varepsilon}$, $\eta_{t,b}^{\xi}$ can be defined so that each dimension of the error processes has idiosyncratic inputs.

In view of the causal representations of the univariate error processes, we have the following causal representations for the high-dimensional nonstationary vector-valued error processes

$$\begin{aligned}\varepsilon_{P,t,n} &= \mathbf{G}_{P,t,n}^\varepsilon(\mathcal{H}_t^\varepsilon) = (G_{P,t,n,i,a}^\varepsilon(\mathcal{H}_{t,a}^\varepsilon))_{i \in [d_X], a \in A_i}, \\ \xi_{P,t,n} &= \mathbf{G}_{P,t,n}^\xi(\mathcal{H}_t^\xi) = (G_{P,t,n,j,b}^\xi(\mathcal{H}_{t,b}^\xi))_{j \in [d_Y], b \in B_j},\end{aligned}$$

where $\mathcal{H}_t^\varepsilon = (\eta_t^\varepsilon, \eta_{t-1}^\varepsilon, \dots)$, $\mathcal{H}_t^\xi = (\eta_t^\xi, \eta_{t-1}^\xi, \dots)$ with $\eta_t^\varepsilon = (\eta_{t,a}^\varepsilon)_{a \in A}$, $\eta_t^\xi = (\eta_{t,b}^\xi)_{b \in B}$ for each $t \in \mathbb{Z}$. Similarly, for each dimension/time-offset tuple $m = (i, j, a, b) \in \mathcal{D}_n$ the error products at time t can be represented as

$$R_{P,t,n,m} = G_{P,t,n,m}^R(\mathcal{H}_{t,m}^R) = G_{P,t,n,i,a}^\varepsilon(\mathcal{H}_{t,a}^\varepsilon) G_{P,t,n,j,b}^\xi(\mathcal{H}_{t,b}^\xi),$$

where $\mathcal{H}_{t,m}^R = (\eta_{t,m}^R, \eta_{t-1,m}^R, \dots)$ with $\eta_{t,m}^R = (\eta_{t,a}^\varepsilon, \eta_{t,b}^\xi)^\top$ for each $t \in \mathbb{Z}$. Also, we have the following representation for the high-dimensional nonstationary \mathbb{R}^{D_n} -valued process of all the products of errors

$$\mathbf{R}_{P,t,n} = \mathbf{G}_{P,t,n}^R(\mathcal{H}_t^R) = (G_{P,t,n,m}^R(\mathcal{H}_{t,m}^R))_{m=(i,j,a,b) \in \mathcal{D}_n}$$

where $\mathcal{H}_t^R = (\eta_t^R, \eta_{t-1}^R, \dots)$ and $\eta_t^R = (\eta_t^\varepsilon, \eta_t^\xi)^\top$ for each $t \in \mathbb{Z}$. Note that for a fixed $P \in \mathcal{P}_n$, $t \in \mathcal{T}_n$, and $n \in \mathbb{N}$ we have that $\mathbf{G}_{P,t,n}^R(\mathcal{H}_s^R)$ is a well-defined high-dimensional random vector for each $s \in \mathbb{Z}$ and $(\mathbf{G}_{P,t,n}^R(\mathcal{H}_s^R))_{s \in \mathbb{Z}}$ is a high-dimensional stationary ergodic \mathbb{R}^{D_n} -valued process.

In view of Assumptions 3.3 and 3.4, for $m = (i, j, a, b) \in \mathcal{D}_n$, we can represent the products of the errors and prediction errors as

$$\begin{aligned}\hat{w}_{P,t,n,m}^{g,\varepsilon} &= G_{P,t,n,m}^{\hat{w}^{g,\varepsilon}}(\mathcal{H}_{t,m}^{\hat{w}^{g,\varepsilon}}) = \hat{w}_{P,t,n,j,b}^{g,\varepsilon} \varepsilon_{P,t,n,i,a}, \\ \hat{w}_{P,t,n,m}^{f,\xi} &= G_{P,t,n,m}^{\hat{w}^{f,\xi}}(\mathcal{H}_{t,m}^{\hat{w}^{f,\xi}}) = \hat{w}_{P,t,n,i,a}^{f,\xi} \xi_{P,t,n,j,b},\end{aligned}$$

with $\mathcal{H}_{t,m}^{\hat{w}^{g,\varepsilon}} = (\mathcal{H}_{t,b}^{\hat{g}}, \mathcal{H}_{t,a}^\varepsilon)$ and $\mathcal{H}_{t,m}^{\hat{w}^{f,\xi}} = (\mathcal{H}_{t,a}^{\hat{f}}, \mathcal{H}_{t,b}^\xi)$. Putting it all together, we have the following causal representation for all dimensions and time-offsets of the products of errors and prediction errors

$$\begin{aligned}\hat{w}_{P,t,n}^{g,\varepsilon} &= \mathbf{G}_{P,t,n}^{\hat{w}^{g,\varepsilon}}(\mathcal{H}_t^{\hat{w}^{g,\varepsilon}}) = (\hat{w}_{P,t,n,m}^{g,\varepsilon})_{m=(i,j,a,b) \in \mathcal{D}_n}, \\ \hat{w}_{P,t,n}^{f,\xi} &= \mathbf{G}_{P,t,n}^{\hat{w}^{f,\xi}}(\mathcal{H}_t^{\hat{w}^{f,\xi}}) = (\hat{w}_{P,t,n,m}^{f,\xi})_{m=(i,j,a,b) \in \mathcal{D}_n},\end{aligned}$$

with $\mathcal{H}_t^{\hat{w}^{g,\varepsilon}} = (\mathcal{H}_t^{\hat{g}}, \mathcal{H}_t^\varepsilon)$ and $\mathcal{H}_t^{\hat{w}^{f,\xi}} = (\mathcal{H}_t^{\hat{f}}, \mathcal{H}_t^\xi)$, where we have suppressed the dependence on n . $\mathbf{G}_{P,t,n}^{\hat{w}^{g,\varepsilon}}(\cdot)$ and $\mathbf{G}_{P,t,n}^{\hat{w}^{f,\xi}}(\cdot)$ are measurable functions such that $\mathbf{G}_{P,t,n}^{\hat{w}^{g,\varepsilon}}(\mathcal{H}_s^{\hat{w}^{g,\varepsilon}})$, $\mathbf{G}_{P,t,n}^{\hat{w}^{f,\xi}}(\mathcal{H}_s^{\hat{w}^{f,\xi}})$ are well-defined high-dimensional random vectors for each $s \in \mathbb{Z}$ and $(\mathbf{G}_{P,t,n}^{\hat{w}^{g,\varepsilon}}(\mathcal{H}_s^{\hat{w}^{g,\varepsilon}}))_{s \in \mathbb{Z}}$, $(\mathbf{G}_{P,t,n}^{\hat{w}^{f,\xi}}(\mathcal{H}_s^{\hat{w}^{f,\xi}}))_{s \in \mathbb{Z}}$ are high-dimensional stationary ergodic processes.

3.4 Assumptions on dependence and nonstationarity

In this subsection, we impose mild assumptions on the rate of decay in temporal dependence and the degree of nonstationarity of the *error processes*. Crucially, these assumptions are stated in a distribution-uniform manner, which is essential for applying the strong Gaussian approximation in Section B. This will be further elaborated upon in Subsection 3.5.

We quantify temporal dependence using the functional dependence measure of Wu [170]. Let $(\tilde{\eta}_{t,a}^\varepsilon, \tilde{\eta}_{t,b}^\xi)_{t \in \mathbb{Z}}$ be an iid copy of $(\eta_{t,a}^\varepsilon, \eta_{t,b}^\xi)_{t \in \mathbb{Z}}$. Denote the set of well-defined tuples of error processes, dimensions, and time-offsets by

$$\mathbb{E} = \{(\varepsilon, i, a) : i \in [d_X], a \in A_i\} \cup \{(\xi, j, b) : j \in [d_Y], b \in B_j\}.$$

For any tuple $(e, l, d) \in \mathbb{E}$ corresponding to a well-defined combination of an error process, dimension, and time-offset, define

$$\tilde{\mathcal{H}}_{t,d,h}^e = (\eta_{t,d}^e, \dots, \eta_{t-h+1,d}^e, \tilde{\eta}_{t-h,d}^e, \eta_{t-h-1,d}^e, \dots)$$

to be $\mathcal{H}_{t,d}^e$ with the input $\eta_{t-h,d}^e$ replaced with the iid copy $\tilde{\eta}_{t-h,d}^e$. Similarly, define $\tilde{\mathcal{H}}_{t,m,h}^R$ as $\mathcal{H}_{t,m}^R$ with the input $\eta_{t-h,m}^R$ replaced with the iid copy $\tilde{\eta}_{t-h,m}^R$ for $m = (i, j, a, b) \in \mathcal{D}_n$, and define $\tilde{\mathcal{H}}_{t,h}^R$ as \mathcal{H}_t^R with the input η_{t-h}^R replaced with the iid copy $\tilde{\eta}_{t-h}^R$. Next, we define the functional dependence measures.

Definition 3.1 (Functional dependence measure). *We define the following measures of temporal dependence for each $n \in \mathbb{N}$, $P \in \mathcal{P}_n$, and $t \in \mathcal{T}_n$. First, define the L^∞ version of the functional dependence measure for the error processes $G_{P,t,n,l,d}^e(\mathcal{H}_{t,d}^e)$ for each $(e, l, d) \in \mathbb{E}$ with $h \in \mathbb{N}_0$ as*

$$\theta_{P,t,n,l,d}^{e,\infty}(h) = \inf\{K \geq 0 : \mathbb{P}_P(|G_{P,t,n,l,d}^e(\mathcal{H}_{t,d}^e) - G_{P,t,n,l,d}^e(\tilde{\mathcal{H}}_{t,d,h}^e)| > K) = 0\}.$$

Second, define the functional dependence measures for the processes of error products $G_{P,t,n,m}^R(\mathcal{H}_{t,m}^R)$ for each $m = (i, j, a, b) \in \mathcal{D}_n$ with $h \in \mathbb{N}_0$, and some $q \geq 1$ as

$$\theta_{P,t,n,m}^R(h, q) = [\mathbb{E}_P(|G_{P,t,n,m}^R(\mathcal{H}_{t,m}^R) - G_{P,t,n,m}^R(\tilde{\mathcal{H}}_{t,m,h}^R)|^q)]^{1/q},$$

and for the vector-valued process $\mathbf{G}_{P,t,n}^R(\mathcal{H}_t^R)$ with $h \in \mathbb{N}_0$, and some $q \geq 1$, $r \geq 1$ as

$$\theta_{P,t,n}^R(h, q, r) = [\mathbb{E}_P(\|\mathbf{G}_{P,t,n}^R(\mathcal{H}_t^R) - \mathbf{G}_{P,t,n}^R(\tilde{\mathcal{H}}_{t,h}^R)\|_r^q)]^{1/q}.$$

For some sequence of collections of distributions $(\mathcal{P}_n)_{n \in \mathbb{N}}$, we make the following assumption about the temporal dependence. We only require the relatively mild assumption that it decays polynomially, rather than geometrically. Note that we will often write the time of the input sequence as 0 when it does not matter due to stationarity.

Assumption 3.5 (Distribution-uniform decay of temporal dependence). *Assume that there exist $\bar{\Theta}^\infty > 0$, $\bar{\beta}^\infty > 1$ such that for all $n \in \mathbb{N}$, $t \in \mathcal{T}_n$, and error processes $(e, l, d) \in \mathbb{E}$, it holds that*

$$\sup_{P \in \mathcal{P}_n} \|G_{P,t,n,l,d}^e(\mathcal{H}_{0,d}^e)\|_{L^\infty(P)} \leq \bar{\Theta}^\infty, \quad \sup_{P \in \mathcal{P}_n} \theta_{P,t,n,l,d}^{e,\infty}(h) \leq \bar{\Theta}^\infty \cdot (h \vee 1)^{-\bar{\beta}^\infty}, \quad h \geq 0.$$

For additional control in terms of the product of errors alone, also assume that there exist $\bar{\Theta}^R > 0$, $\bar{\beta}^R > 3$, $\bar{q}^R > 4$, such that for all $n \in \mathbb{N}$, $t \in \mathcal{T}_n$, $m = (i, j, a, b) \in \mathcal{D}_n$, it holds that

$$\sup_{P \in \mathcal{P}_n} [\mathbb{E}_P(|G_{P,t,n,m}^R(\mathcal{H}_{0,m}^R)|^{\bar{q}^R})]^{1/\bar{q}^R} \leq \bar{\Theta}^R, \quad \sup_{P \in \mathcal{P}_n} \theta_{P,t,n,m}^R(h, \bar{q}^R) \leq \bar{\Theta}^R \cdot (h \vee 1)^{-\bar{\beta}^R}, \quad h \geq 0.$$

A few remarks are in order. First, the constants in Assumption 3.5 do not depend on n . Second, the assumptions on the individual error processes can be weakened; see Subsection A.7 for more discussion. Third, for all $n \in \mathbb{N}$, $t \in \mathcal{T}_n$, by Jensen's inequality we have

$$\sup_{P \in \mathcal{P}_n} [\mathbb{E}_P(\|\mathbf{G}_{P,t,n}^R(\mathcal{H}_0^R)\|_2^{\bar{q}^R})]^{1/\bar{q}^R} \leq D_n^{\frac{1}{2}} \bar{\Theta}^R, \quad \sup_{P \in \mathcal{P}_n} \theta_{P,t,n}^R(h, \bar{q}^R, 2) \leq D_n^{\frac{1}{2}} \bar{\Theta}^R \cdot (h \vee 1)^{-\bar{\beta}^R}, \quad h \geq 0.$$

Next, for some sequence of collections of distributions $(\mathcal{P}_n)_{n \in \mathbb{N}}$, we make the following assumption to control the nonstationarity of the process of error products.

Assumption 3.6 (Distribution-uniform total variation condition for nonstationarity). *Recall $\bar{\Theta}^R > 0$ from Assumption 3.5. Assume that for each $n \in \mathbb{N}$, there exists a constant $\bar{\Gamma}_n^R \geq 1$ such that*

$$\sup_{P \in \mathcal{P}_n} \left(\sum_{t=\mathbb{T}_n+1}^{\mathbb{T}_n^+} (\mathbb{E}_P \|\mathbf{G}_{P,t,n}^R(\mathcal{H}_0^R) - \mathbf{G}_{P,t-1,n}^R(\mathcal{H}_0^R)\|_2^2)^{1/2} \right) \leq \bar{\Theta}^R \bar{\Gamma}_n^R.$$

3.5 Theoretical result for dGCM

In this subsection, we present the theoretical result that justifies the bootstrap procedure described in Algorithm 1. This result relies on time-varying nonlinear regression for nonstationary processes and the distribution-uniform strong Gaussian approximation from Section B applied to the process of error products. The approximating nonstationary Gaussian process has a time-varying covariance structure, which is explicitly characterized by *local long-run* covariance matrices.

Definition 3.2 (Local long-run covariance matrices of the process of error products). *For each $P \in \mathcal{P}_n$, $t \in \mathcal{T}_n$, $n \in \mathbb{N}$, define the local long-run covariance matrix $\Sigma_{P,t,n}^R \in \mathbb{R}^{D_n \times D_n}$ of the \mathbb{R}^{D_n} -valued stationary process $(\mathbf{G}_{P,t,n}^R(\mathcal{H}_t^R))_{t \in \mathbb{Z}}$ by*

$$\Sigma_{P,t,n}^R = \sum_{h \in \mathbb{Z}} \text{Cov}_P(\mathbf{G}_{P,t,n}^R(\mathcal{H}_0^R), \mathbf{G}_{P,t,n}^R(\mathcal{H}_h^R)).$$

In view of the Gaussian approximation theory developed in Mies and Steland [109], we only require an estimator of the *cumulative* covariance matrices of the error products

$$Q_{P,t,n}^{\mathbf{R}} = \sum_{s=\mathbb{T}_n^-}^t \Sigma_{P,s,n}^{\mathbf{R}},$$

rather than the local long-run covariance matrices at each time individually. This is critical for the practical applicability of our method, as estimating individual local long-run covariance matrices can be extremely challenging in practice. Specifically, we use the estimator

$$\hat{Q}_{t,n}^{\mathbf{R}} = \sum_{s=L_n+\mathbb{T}_n^--1}^t \frac{1}{L_n} \left(\sum_{r=s-L_n+1}^s \hat{\mathbf{R}}_{r,n} \right)^{\otimes 2}, \quad (12)$$

for some lag-window size $L_n \in \mathbb{N}$. We discuss how to select L_n in practice using the minimum volatility method in Subsection 5.1. Going forward, denote $Q_{P,n}^{\mathbf{R}} = (Q_{P,t,n}^{\mathbf{R}})_{t \in \mathcal{T}_{n,L}}$ and $\hat{Q}_n^{\mathbf{R}} = (\hat{Q}_{t,n}^{\mathbf{R}})_{t \in \mathcal{T}_{n,L}}$, where

$$\mathcal{T}_{n,L} = \{L_n + \mathbb{T}_n^- - 1, \dots, \mathbb{T}_n^+ - 1, \mathbb{T}_n^+\}.$$

To account for the estimation errors for the time-varying regression functions and the cumulative covariance matrices, as well as the error for the Gaussian approximation, we introduce offsets $\tau_n \rightarrow 0$, $\nu_n \rightarrow 0$ so that $\tau_n = o(\log^{-(1+\delta)}(T_n))$ for some $\delta > 0$ and

$$\nu_n \gg \log(T_n) D_n \left[\left(\frac{D_n}{T_n} \right)^{2\xi(\bar{q}^R, \bar{\beta}^R)} + \tau_n^{-2} (\varphi_{n,1} + \varphi_{n,2}) \right], \quad (13)$$

where

$$\varphi_{n,1} = T_n^{-\frac{1}{2}} (\bar{\Gamma}_n^R)^{\frac{1}{2}} L_n^{\frac{1}{4}} + T_n^{-\frac{1}{4}} D_n^{\frac{1}{4}} L_n^{\frac{1}{4}} + L_n^{-\frac{1}{2}} + L_n^{1-\frac{\bar{\beta}^R}{2}} + T_n^{-1},$$

comes from the covariance estimation error,

$$\varphi_{n,2} = \tau_n^{\frac{7}{2}} D_n^{-\frac{5}{4}} + \tau_n^7 D_n^{-\frac{5}{2}},$$

comes from the time-varying regression function estimation errors. Also, the lag-window size L_n from (12) must satisfy $L_n \asymp T_n^\zeta$ for some $\zeta \in (0, \frac{1}{2})$ so that $\tau_n^{-6} D_n^2 L_n^{-1} = o(1)$ and $\bar{\Gamma}_n^R T_n^{-1} D_n^2 \tau_n^{-6} L_n^{\frac{1}{2}} = o(1)$, where $\bar{\Gamma}_n^R$ is from Assumption 3.6. We see that the offsets depend on the number of observations T_n from Subsection 2.1, the intrinsic dimensionality D_n from Subsection 2.1, the degree of nonstationarity $\bar{\Gamma}_n^R$ from Assumption 3.6, and the lag-window parameter L_n from (12). $\xi(\bar{q}^R, \bar{\beta}^R)$ is a rate defined in Section B that depends on the constants $\bar{\beta}^R, \bar{q}^R$ from Assumption 3.5.

The following result establishes the validity of our bootstrap-based testing procedure described in Algorithm 1, provided that the previously stated assumptions hold and the prediction errors

$$\begin{aligned} \hat{w}_{P,t,n,i,a}^f &= f_{P,t,n,i,a}(\mathbf{Z}_{t,n}) - \hat{f}_{t,n,i,a}(\mathbf{Z}_{t,n}), \\ \hat{w}_{P,t,n,j,b}^g &= g_{P,t,n,j,b}(\mathbf{Z}_{t,n}) - \hat{g}_{t,n,j,b}(\mathbf{Z}_{t,n}), \end{aligned}$$

converge to zero sufficiently fast, in some sense. If it were known, we could correctly calibrate our test with the (random) quantile function \hat{q} of $S_{n,p}(\check{\mathbf{R}}_n)$, where $\check{\mathbf{R}}_n = (\check{\mathbf{R}}_{t,n})_{t \in \mathcal{T}_{n,L}}$ and $\check{\mathbf{R}}_{t,n} \sim \mathcal{N}(0, \hat{\Sigma}_{t,n}^{\mathbf{R}})$ for all $t \in \mathcal{T}_{n,L}$. In practice, \hat{q} is numerically approximated by conducting a large number of Monte Carlo simulations, and we use \hat{q}^{boot} from Algorithm 1 in its place.

Theorem 3.1. *Suppose that Assumptions 3.1, 3.2, 3.3, 3.4, 3.5, 3.6 related to the temporal dependence and nonstationarity of the processes all hold for the sequence of collections of distributions $(\mathcal{P}_{0,n}^*)_{n \in \mathbb{N}}$, where $\mathcal{P}_{0,n}^* \subset \mathcal{P}_{0,n}^{\text{CI}}$ for each $n \in \mathbb{N}$. Further, suppose that*

$$\begin{aligned} \sup_{P \in \mathcal{P}_{0,n}^*} \max_{(i,j,a,b) \in \mathcal{D}_n} \max_{t \in \mathcal{T}_n} \mathbb{E}_P \left(\left| \hat{w}_{P,t,n,i,a}^f \right|^2 \right)^{\frac{1}{2}} \mathbb{E}_P \left(\left| \hat{w}_{P,t,n,j,b}^g \right|^2 \right)^{\frac{1}{2}} &= o(T_n^{-\frac{1}{2}} \tau_n^7 D_n^{-\frac{3}{2}}), \\ \sup_{P \in \mathcal{P}_{0,n}^*} \max_{i \in [d_X], a \in A_i} \max_{t \in \mathcal{T}_n} \mathbb{E}_P \left(\left| \hat{w}_{P,t,n,i,a}^f \right|^2 \right)^{\frac{1}{2}} &= o(\tau_n^7 D_n^{-\frac{5}{2}}), \\ \sup_{P \in \mathcal{P}_{0,n}^*} \max_{j \in [d_Y], b \in B_j} \max_{t \in \mathcal{T}_n} \mathbb{E}_P \left(\left| \hat{w}_{P,t,n,j,b}^g \right|^2 \right)^{\frac{1}{2}} &= o(\tau_n^7 D_n^{-\frac{5}{2}}). \end{aligned}$$

If the offsets $\tau_n \rightarrow 0$ and $\nu_n \rightarrow 0$ are chosen such that condition (13) holds, then we have

$$\limsup_{n \rightarrow \infty} \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P \left(S_{n,p}(\hat{\mathbf{R}}_n) > \hat{q}_{1-\alpha+\nu_n} + \tau_n \right) \leq \alpha.$$

The previous result demonstrates that the dGCM test possesses a property known as *rate double robustness*. This property means that we only require modest convergence rates for the *products* of the prediction errors, rather than for each prediction error individually. This feature of the dGCM test can be especially useful in the contexts of causal discovery for time-lagged effects and variable selection in time series forecasting. In these applications, a faster convergence rate of a nowcasting model can compensate for a slower convergence rate of a forecasting model, or vice versa.

4 dGCM with Sieve Time-Varying Regression (Sieve-dGCM)

The purpose of this section is to demonstrate that the convergence rates required by Theorem 3.1 for estimating the time-varying regression functions can be achieved. To show this, we consider an instantiation of the dynamic generalized covariance measure (dGCM) test based on the sieve time-varying nonlinear regression estimator from Ding and Zhou [48]. We refer to this instantiation of the dGCM test as the Sieve-dGCM test.

We prove that, under mild assumptions about the temporal dependence and nonstationarity of the processes, the sieve estimator achieves the required convergence rates and the Sieve-dGCM test has asymptotic Type-I error control. In Section 5, we study the finite sample performance of the Sieve-dGCM test. Along the way, in Subsection 5.1, we introduce a novel cross-validation scheme which we use for selecting the parameters of the sieve estimator.

In this section, we use the framework of *locally stationary processes* [36, 190, 37, 40]. This is a well-studied class of nonstationary processes that fits within the general triangular array framework for nonstationary processes from Section 3. We note that there are several other time-varying regression estimators for locally stationary processes; see Zhang and Wu [186], Yousuf and Ng [177], and Chen, Smetanina, and Wu [27] and the references therein. We discuss extensions to *piecewise* locally stationary processes in Subsection A.6.

4.1 Setting and notation

We follow Dahlhaus [36] in rescaling time to the unit interval $t/n \in [0, 1]$, so that *infill asymptotics* can be used to study nonstationary processes. In this setting, the sample size n no longer corresponds to getting information about the future. Instead, as n increases we get more observations about each *local structure* of the nonstationary process. Zhou and Wu [190] introduced the framework for representing locally stationary processes as nonlinear functions of iid inputs as in Wu [170].

We use the same notation as Subsection 2.1, with the only difference being that we fix the number of dimensions d_Z and time-offsets C_k for each dimension $k \in [d_Z]$. We still allow the number of dimensions $d_X = d_{X,n}$, $d_Y = d_{Y,n}$ and time-offsets A_i , B_j for each $i \in [d_X]$, $j \in [d_Y]$ to grow with n . Define A , B , C as the collection of all time-offsets as in Subsection 2.1, where $A = A_n$, $B = B_n$ and C is fixed. We emphasize that there is no inherent necessity for fixing the number of dimensions d_Z and time-offsets C_k . Our reason for doing this is because we want to leverage the *existing* theoretical results for the sieve estimator from Ding and Zhou [48]. Future investigations can study the performance of the sieve estimator in the high-dimensional setting, so that we can allow the number of dimensions d_Z and time-offsets C_k to grow with n .

We will still use the notation \mathcal{T}_n for the subset of original times in which *all* time-offsets of each dimension of $X_{t,n}$, $Y_{t,n}$, and $Z_{t,n}$ are actually observed,

$$\mathcal{T}_n = \{1 - \min(a_{\min}, b_{\min}, c_{\min}), n - \max(a_{\max}, b_{\max}, c_{\max})\} \subseteq \{1, \dots, n\}.$$

Also, we will still denote $T_n = |\mathcal{T}_n|$, $\mathbb{T}_n^- = \min(\mathcal{T}_n)$, and $\mathbb{T}_n^+ = \max(\mathcal{T}_n)$. Similarly, denote the corresponding interval of rescaled times in which all time-offsets are well-defined by

$$\mathcal{U}_n = \left[\frac{1}{n} - \frac{\min(a_{\min}, b_{\min}, c_{\min})}{n}, 1 - \frac{\max(a_{\max}, b_{\max}, c_{\max})}{n} \right] \subset [0, 1],$$

and denote $\mathbb{U}_n^- = \min(\mathcal{U}_n)$, and $\mathbb{U}_n^+ = \max(\mathcal{U}_n)$.

Recall the index set containing the dimensions and time-offsets of interest

$$\mathcal{D}_n \subseteq \{(i, j, a, b) : i \in [d_X], j \in [d_Y], a \in A_i, b \in B_j\},$$

where A_i, B_j are the time-offsets for dimensions $i \in [d_X], j \in [d_Y]$. Again, we will often refer to the dimension/time-offset tuple by $m = (i, j, a, b) \in \mathcal{D}_n$ to lighten the notation. Denote cardinality $D_n = |\mathcal{D}_n|$ which may grow with n .

4.2 Locally stationary observed processes

Next, we introduce the causal representation of locally stationary processes, which is most similar to Ding and Zhou [48] and Example 3 in Mies and Steland [109]. This representation is different than the previous causal representation from Assumption 3.1, because we now assume that the nonlinear stochastic system is well-defined for *all* rescaled times. For the following assumption, let

$$\mathcal{H}_t^X = (\eta_t^X, \eta_{t-1}^X, \dots), \mathcal{H}_t^Y = (\eta_t^Y, \eta_{t-1}^Y, \dots), \mathcal{H}_t^Z = (\eta_t^Z, \eta_{t-1}^Z, \dots),$$

where $(\eta_t^X, \eta_t^Y, \eta_t^Z)_{t \in \mathbb{Z}}$ is an iid sequence of random vectors. Denote the dimensions of $\eta_t^X = \eta_{t,n}^X$, $\eta_t^Y = \eta_{t,n}^Y$, $\eta_t^Z = \eta_{t,n}^Z$ respectively by $d_X^\eta = d_{X,n}^\eta$, $d_Y^\eta = d_{Y,n}^\eta$, $d_Z^\eta = d_{Z,n}^\eta$, which can change with n .

Assumption 4.1 (Causal representations of the observed processes). *Assume that we can represent each dimension of each of the observed processes as the output of an evolving nonlinear system that was given a sequence of iid inputs:*

$$X_{t,n,i} = \tilde{G}_{n,i}^X(t/n, \mathcal{H}_t^X), Y_{t,n,j} = \tilde{G}_{n,j}^Y(t/n, \mathcal{H}_t^Y), Z_{t,n,k} = \tilde{G}_{n,k}^Z(t/n, \mathcal{H}_t^Z),$$

where the systems are defined for all $u \in [0, 1]$ by

$$\tilde{X}_{t,n,i}(u) = \tilde{G}_{n,i}^X(u, \mathcal{H}_t^X), \tilde{Y}_{t,n,j}(u) = \tilde{G}_{n,j}^Y(u, \mathcal{H}_t^Y), \tilde{Z}_{t,n,k}(u) = \tilde{G}_{n,k}^Z(u, \mathcal{H}_t^Z),$$

so that we have $X_{t,n,i} = \tilde{X}_{t,n,i}(t/n)$, $Y_{t,n,j} = \tilde{Y}_{t,n,j}(t/n)$, $Z_{t,n,k} = \tilde{Z}_{t,n,k}(t/n)$.

For each $n \in \mathbb{N}$, $(i, j, a, b) \in \mathcal{D}_n$, $t \in \mathcal{T}_n$, we assume that $\tilde{G}_{n,i}^X(u, \cdot)$, $\tilde{G}_{n,j}^Y(u, \cdot)$, $\tilde{G}_{n,k}^Z(u, \cdot)$ are measurable functions from $(\mathbb{R}^{d_X^\eta})^\infty$, $(\mathbb{R}^{d_Y^\eta})^\infty$, $(\mathbb{R}^{d_Z^\eta})^\infty$, respectively, to \mathbb{R} — where we endow $(\mathbb{R}^{d_X^\eta})^\infty$, $(\mathbb{R}^{d_Y^\eta})^\infty$, $(\mathbb{R}^{d_Z^\eta})^\infty$ with the σ -algebra generated by all finite projections — such that $\tilde{G}_{n,i}^X(u, \mathcal{H}_s^X)$, $\tilde{G}_{n,j}^Y(u, \mathcal{H}_s^Y)$, $\tilde{G}_{n,k}^Z(u, \mathcal{H}_s^Z)$ are each well-defined random variables for each $s \in \mathbb{Z}$ and $(\tilde{G}_{n,i}^X(u, \mathcal{H}_s^X))_{s \in \mathbb{Z}}$, $(\tilde{G}_{n,j}^Y(u, \mathcal{H}_s^Y))_{s \in \mathbb{Z}}$, $(\tilde{G}_{n,k}^Z(u, \mathcal{H}_s^Z))_{s \in \mathbb{Z}}$ are each stationary ergodic processes.

As in Subsection 3.1, we have not defined the input sequences for the observed processes separately for each dimension. However, without loss of generality, we can define the measurable functions $\tilde{G}_{n,i}^X(u, \cdot)$, $\tilde{G}_{n,j}^Y(u, \cdot)$, $\tilde{G}_{n,k}^Z(u, \cdot)$ and the inputs η_t^X , η_t^Y , η_t^Z so that each dimension of the observed processes can have idiosyncratic inputs.

In light of Assumption 4.1, we have the following causal representations for all dimensions with no time-offsets

$$\begin{aligned} \tilde{X}_{t,n}(u) &= \tilde{G}_n^X(u, \mathcal{H}_t^X) = (\tilde{G}_{n,i}^X(u, \mathcal{H}_t^X))_{i \in [d_X]}, \\ \tilde{Y}_{t,n}(u) &= \tilde{G}_n^Y(u, \mathcal{H}_t^Y) = (\tilde{G}_{n,j}^Y(u, \mathcal{H}_t^Y))_{j \in [d_Y]}, \\ \tilde{Z}_{t,n}(u) &= \tilde{G}_n^Z(u, \mathcal{H}_t^Z) = (\tilde{G}_{n,k}^Z(u, \mathcal{H}_t^Z))_{k \in [d_Z]}, \end{aligned}$$

so that we have $X_{t,n} = \tilde{X}_{t,n}(t/n)$, $Y_{t,n} = \tilde{Y}_{t,n}(t/n)$, $Z_{t,n} = \tilde{Z}_{t,n}(t/n)$. For each $n \in \mathbb{N}$, we have causal representations for dimensions $i \in [d_X]$, $j \in [d_Y]$, $k \in [d_Z]$ with time-offsets $a \in A_i$, $b \in B_j$, $c \in C_k$

$$\begin{aligned} \tilde{X}_{t,n,i,a}(u) &= \tilde{G}_{n,i,a}^X(u, \mathcal{H}_{t,a}^X) = \tilde{G}_{n,i}^X(u + a/n, \mathcal{H}_{t+a}^X), \\ \tilde{Y}_{t,n,j,b}(u) &= \tilde{G}_{n,j,b}^Y(u, \mathcal{H}_{t,b}^Y) = \tilde{G}_{n,j}^Y(u + b/n, \mathcal{H}_{t+b}^Y), \\ \tilde{Z}_{t,n,k,c}(u) &= \tilde{G}_{n,k,c}^Z(u, \mathcal{H}_{t,c}^Z) = \tilde{G}_{n,k}^Z(u + c/n, \mathcal{H}_{t+c}^Z), \end{aligned}$$

where $\mathcal{H}_{t,a}^X = (\eta_{t+a}^X, \eta_{t-1+a}^X, \dots)$, $\mathcal{H}_{t,b}^Y = (\eta_{t+b}^Y, \eta_{t-1+b}^Y, \dots)$, and $\mathcal{H}_{t,c}^Z = (\eta_{t+c}^Z, \eta_{t-1+c}^Z, \dots)$, so that we have $X_{t,n,i,a} = \tilde{X}_{t,n,i,a}(t/n)$, $Y_{t,n,j,b} = \tilde{Y}_{t,n,j,b}(t/n)$, $Z_{t,n,k,c} = \tilde{Z}_{t,n,k,c}(t/n)$ for each dimension of the

observed sequence with time-offset. We can then write the causal representation of the vectors with all dimensions and time-offsets as

$$\begin{aligned}\tilde{\mathbf{X}}_{t,n}(u) &= \tilde{\mathbf{G}}_n^{\mathbf{X}}(u, \mathcal{H}_t^{\mathbf{X}}) = (\tilde{G}_{n,i,a}^{\mathbf{X}}(u, \mathcal{H}_{t,a}^{\mathbf{X}}))_{i \in [d_{\mathbf{X}}], a \in A_i}, \\ \tilde{\mathbf{Y}}_{t,n}(u) &= \tilde{\mathbf{G}}_n^{\mathbf{Y}}(u, \mathcal{H}_t^{\mathbf{Y}}) = (\tilde{G}_{n,j,b}^{\mathbf{Y}}(u, \mathcal{H}_{t,b}^{\mathbf{Y}}))_{j \in [d_{\mathbf{Y}}], b \in B_j}, \\ \tilde{\mathbf{Z}}_{t,n}(u) &= \tilde{\mathbf{G}}_n^{\mathbf{Z}}(u, \mathcal{H}_t^{\mathbf{Z}}) = (\tilde{G}_{n,k,c}^{\mathbf{Z}}(u, \mathcal{H}_{t,c}^{\mathbf{Z}}))_{k \in [d_{\mathbf{Z}}], c \in C_k},\end{aligned}$$

where $\mathcal{H}_t^{\mathbf{X}} = (\eta_t^{\mathbf{X}}, \eta_{t-1}^{\mathbf{X}}, \dots)$, $\mathcal{H}_t^{\mathbf{Y}} = (\eta_t^{\mathbf{Y}}, \eta_{t-1}^{\mathbf{Y}}, \dots)$, $\mathcal{H}_t^{\mathbf{Z}} = (\eta_t^{\mathbf{Z}}, \eta_{t-1}^{\mathbf{Z}}, \dots)$, and $\eta_t^{\mathbf{X}} = \eta_{t+a_{\max}}^{\mathbf{X}}$, $\eta_t^{\mathbf{Y}} = \eta_{t+b_{\max}}^{\mathbf{Y}}$, $\eta_t^{\mathbf{Z}} = \eta_{t+c_{\max}}^{\mathbf{Z}}$, so that we have $\mathbf{X}_{t,n} = \tilde{\mathbf{X}}_{t,n}(t/n)$, $\mathbf{Y}_{t,n} = \tilde{\mathbf{Y}}_{t,n}(t/n)$, $\mathbf{Z}_{t,n} = \tilde{\mathbf{Z}}_{t,n}(t/n)$ for the observed sequence including all dimensions and time-offsets.

Let Ω be a sample space, \mathcal{B} the Borel sigma-algebra, and (Ω, \mathcal{B}) a measurable space. For fixed $n \in \mathbb{N}$, let (Ω, \mathcal{B}) be equipped with a family of probability measures $(\mathbb{P}_P)_{P \in \mathcal{P}_n}$ so that the joint distribution of the nonlinear stochastic systems

$$(\tilde{G}_n^{\mathbf{X}}(u, \mathcal{H}_t^{\mathbf{X}}))_{u \in [0,1], t \in \mathbb{Z}}, (\tilde{G}_n^{\mathbf{Y}}(u, \mathcal{H}_t^{\mathbf{Y}}))_{u \in [0,1], t \in \mathbb{Z}}, (\tilde{G}_n^{\mathbf{Z}}(u, \mathcal{H}_t^{\mathbf{Z}}))_{u \in [0,1], t \in \mathbb{Z}}$$

under \mathbb{P}_P is $P \in \mathcal{P}_n$, where the collection of distributions \mathcal{P}_n can change with n . The family of probability measures $(\mathbb{P}_P)_{P \in \mathcal{P}_n}$ is defined with respect to the same measurable space (Ω, \mathcal{B}) , but need not have the same dominating measure.

We use the same null hypotheses of conditional independence as those in Subsection 2.2. Again, for each $n \in \mathbb{N}$, we denote the collection of distributions such that the null hypothesis is true by $\mathcal{P}_{0,n}^{\text{CI}}$. In the locally stationary setting, the null hypothesis

$$X_{t,n,i,a} \perp\!\!\!\perp Y_{t,n,j,b} \mid \mathbf{Z}_{t,n} \text{ for all } t \in \mathcal{T}_n, \text{ for all } (i, j, a, b) \in \mathcal{D}_n, \quad (14)$$

can be written equivalently as

$$\tilde{X}_{t,n,i,a}(t/n) \perp\!\!\!\perp \tilde{Y}_{t,n,j,b}(t/n) \mid \tilde{\mathbf{Z}}_{t,n}(t/n) \text{ for all } t \in \mathcal{T}_n, \text{ for all } (i, j, a, b) \in \mathcal{D}_n,$$

where \mathcal{D}_n only contains a single dimension/time-offset tuple in the univariate setting.

We will state more assumptions in the next several subsections for a generic sequence of collections of distributions $(\mathcal{P}_n)_{n \in \mathbb{N}}$. In Theorem 4.1, we will assume that these conditions hold for the sequence of collections of distributions $(\mathcal{P}_{0,n}^*)_{n \in \mathbb{N}}$, where $\mathcal{P}_{0,n}^* \subset \mathcal{P}_{0,n}^{\text{CI}}$ for each $n \in \mathbb{N}$. Note that we make stronger assumptions in this section than in Section 3 to ensure that the sieve estimators satisfy the convergence rate requirements of Theorem 3.1.

4.3 Sieve time-varying nonlinear regression estimator

For a given sample size $n \in \mathbb{N}$, distribution $P \in \mathcal{P}_n$, time $t \in \mathcal{T}_n$, and dimension/time-offset tuple $(i, j, a, b) \in \mathcal{D}_n$, we consider the time-varying nonlinear regression model

$$X_{t,n,i,a} = f_{P,n,i,a}(t/n, \mathbf{Z}_{t,n}) + \varepsilon_{P,t,n,i,a}, \quad Y_{t,n,j,b} = g_{P,n,j,b}(t/n, \mathbf{Z}_{t,n}) + \xi_{P,t,n,j,b},$$

where $f_{P,n,i,a}(u, \mathbf{z})$ and $g_{P,n,j,b}(u, \mathbf{z})$ are smooth functions of rescaled time u and covariate values \mathbf{z} with $f_{P,n,i,a}(t/n, \mathbf{z}) = \mathbb{E}_P(X_{t,n,i,a} \mid \mathbf{Z}_{t,n} = \mathbf{z})$ and $g_{P,n,j,b}(t/n, \mathbf{z}) = \mathbb{E}_P(Y_{t,n,j,b} \mid \mathbf{Z}_{t,n} = \mathbf{z})$. We emphasize that the functions $f_{P,n,i,a}(u, \mathbf{z})$ and $g_{P,n,j,b}(u, \mathbf{z})$ depend on rescaled time u rather than “real time” t , as in the literature on nonparametric regression for locally stationary processes [161, 186, 177, 27, 48]. For $m = (i, j, a, b) \in \mathcal{D}_n$, denote the error products at time t by

$$R_{P,t,n,m} = \varepsilon_{P,t,n,i,a} \xi_{P,t,n,j,b},$$

and the corresponding residual products by

$$\hat{R}_{t,n,m} = \hat{\varepsilon}_{t,n,i,a} \hat{\xi}_{t,n,j,b},$$

where $\hat{\varepsilon}_{t,n,i,a} = X_{t,n,i,a} - \hat{f}_{t,n,i,a}(t/n, \mathbf{Z}_{t,n})$ and $\hat{\xi}_{t,n,j,b} = Y_{t,n,j,b} - \hat{g}_{t,n,j,b}(t/n, \mathbf{Z}_{t,n})$.

The estimates $\hat{f}_{t,n,i,a}$ and $\hat{g}_{t,n,j,b}$ of the functions $f_{P,n,i,a}$ and $g_{P,n,j,b}$ are formed by regressing $(X_{t,n,i,a})_{t \in \mathcal{T}_n}$ on $(\mathbf{Z}_{t,n})_{t \in \mathcal{T}_n}$ and $(Y_{t,n,j,b})_{t \in \mathcal{T}_n}$ on $(\mathbf{Z}_{t,n})_{t \in \mathcal{T}_n}$, respectively, using the time-varying nonlinear sieve regression estimator introduced below. The subscript t in $\hat{f}_{t,n,i,a}$ and $\hat{g}_{t,n,j,b}$ is to indicate that

we allow for *sequential estimation*, which will be discussed in Remark 4.1. Let $\hat{\mathbf{R}}_{t,n} = (\hat{R}_{t,n,m})_{m \in \mathcal{D}_n}$ be the high-dimensional vector process containing the residual products for all dimension/time-offset combinations in \mathcal{D}_n . The observed processes X, Y, Z and error processes ϵ, ξ can all be locally stationary processes; see Subsection 4.5 for the details.

For some sequence of collections of distributions $(\mathcal{P}_n)_{n \in \mathbb{N}}$, we make the following assumption.

Assumption 4.2 (Additive form and regularity). *For each sample size $n \in \mathbb{N}$, distribution $P \in \mathcal{P}_n$, rescaled time $u \in \mathcal{U}_n$, and dimension/time-offset tuple $(i, j, a, b) \in \mathcal{D}_n$, assume that*

$$\begin{aligned} f_{P,n,i,a}(u, \mathbf{z}) &= \sum_{k=1}^{d_Z} \sum_{c=1}^{C_k} f_{P,n,i,a,k,c}(u, z_{k,c}), \\ g_{P,n,j,b}(u, \mathbf{z}) &= \sum_{k=1}^{d_Z} \sum_{c=1}^{C_k} g_{P,n,j,b,k,c}(u, z_{k,c}), \end{aligned}$$

where $f_{P,n,i,a,k,c} : \mathcal{U}_n \times \mathbb{R} \rightarrow \mathbb{R}$ and $g_{P,n,j,b,k,c} : \mathcal{U}_n \times \mathbb{R} \rightarrow \mathbb{R}$ are time-varying partial response functions, so that we have

$$\begin{aligned} \mathbb{E}_P(X_{t,n,i,a} | \mathbf{Z}_{t,n} = \mathbf{z}) &= \sum_{k=1}^{d_Z} \sum_{c=1}^{C_k} f_{P,n,i,a,k,c}(t/n, z_{k,c}), \\ \mathbb{E}_P(Y_{t,n,j,b} | \mathbf{Z}_{t,n} = \mathbf{z}) &= \sum_{k=1}^{d_Z} \sum_{c=1}^{C_k} g_{P,n,j,b,k,c}(t/n, z_{k,c}), \end{aligned}$$

for each time $t \in \mathcal{T}_n$.

Further, assume for all $n \in \mathbb{N}$, $i \in [d_X]$, $a \in A_i$, $j \in [d_Y]$, $b \in B_j$, $k \in [d_Z]$, $c \in C_k$, $u \in \mathcal{U}_n$, there exists some $q \geq 2$ such that

$$\begin{aligned} \sup_{P \in \mathcal{P}_n} \mathbb{E}_P(|f_{P,n,i,a,k,c}(u, Z_{t,n,k,c})|^q) &< \infty, \\ \sup_{P \in \mathcal{P}_n} \mathbb{E}_P(|g_{P,n,j,b,k,c}(u, Z_{t,n,k,c})|^q) &< \infty. \end{aligned}$$

To fix ideas, we use the algebraic mapping $h : [-1, 1] \rightarrow \mathbb{R}$ from Example 3.1 in Ding and Zhou [48] with positive scaling factor $s = 1$,

$$h(\tilde{z}) = \begin{cases} -\infty, & \tilde{z} = -1, \\ \frac{\tilde{z}}{\sqrt{1-\tilde{z}^2}}, & \tilde{z} \in (-1, 1), \\ \infty, & \tilde{z} = 1. \end{cases}$$

See the discussion preceding Definition 3.1 in Ding and Zhou [48] for additional details. For some sequence of collections of distributions $(\mathcal{P}_n)_{n \in \mathbb{N}}$, for each $n \in \mathbb{N}$, $i \in [d_X]$, $a \in A_i$, $j \in [d_Y]$, $b \in B_j$, $k \in [d_Z]$, $c \in C_k$, and $P \in \mathcal{P}_n$, we relate the time-varying partial response functions $f_{P,n,i,a,k,c} : \mathcal{U}_n \times \mathbb{R} \rightarrow \mathbb{R}$ and $g_{P,n,j,b,k,c} : \mathcal{U}_n \times \mathbb{R} \rightarrow \mathbb{R}$ to $\tilde{f}_{P,n,i,a,k,c} : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$ and $\tilde{g}_{P,n,j,b,k,c} : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$, respectively, where

$$\begin{aligned} \tilde{f}_{P,n,i,a,k,c}(u^*, z^*) &= f_{P,n,i,a,k,c}(\mathbb{U}_n^- + u^*(\mathbb{U}_n^+ - \mathbb{U}_n^-), h(2z^* - 1)), \\ \tilde{g}_{P,n,j,b,k,c}(u^*, z^*) &= g_{P,n,j,b,k,c}(\mathbb{U}_n^- + u^*(\mathbb{U}_n^+ - \mathbb{U}_n^-), h(2z^* - 1)), \end{aligned}$$

with $\mathbb{U}_n^- = \min(\mathcal{U}_n)$ and $\mathbb{U}_n^+ = \max(\mathcal{U}_n)$.

For some sequence of collections of distributions $(\mathcal{P}_n)_{n \in \mathbb{N}}$, we make the following assumption.

Assumption 4.3 (Smoothness). *For each $n \in \mathbb{N}$, $i \in [d_X]$, $a \in A_i$, $j \in [d_Y]$, $b \in B_j$, $k \in [d_Z]$, $c \in C_k$, and $P \in \mathcal{P}_n$, assume that for each fixed $u^* \in [0, 1]$ we have*

$$\tilde{f}_{P,n,i,a,k,c}(u^*, \cdot) \in C^\infty([0, 1]), \quad \tilde{g}_{P,n,j,b,k,c}(u^*, \cdot) \in C^\infty([0, 1]),$$

and for each fixed $z^* \in [0, 1]$ we have

$$\tilde{f}_{P,n,i,a,k,c}(\cdot, z^*) \in C^\infty([0, 1]), \quad \tilde{g}_{P,n,j,b,k,c}(\cdot, z^*) \in C^\infty([0, 1]),$$

where $C^\infty([0, 1])$ denotes the space of functions on $[0, 1]$ that are infinitely differentiable.

If Assumption 4.3 holds, then by Theorem 3.1 of Ding and Zhou [48] we can approximate the time-varying partial response functions by

$$\begin{aligned} f_{P,n,i,a,k,c}(u, z) &\approx \sum_{\ell_1=1}^{\tilde{c}_n} \sum_{\ell_2=1}^{\tilde{d}_n} \beta_{P,n,i,a,k,c,\ell_1,\ell_2}^f b_{\ell_1,\ell_2}(u, z), \\ g_{P,n,j,b,k,c}(u, z) &\approx \sum_{\ell_1=1}^{\tilde{c}_n} \sum_{\ell_2=1}^{\tilde{d}_n} \beta_{P,n,j,b,k,c,\ell_1,\ell_2}^g b_{\ell_1,\ell_2}(u, z), \end{aligned}$$

where $\{b_{\ell_1,\ell_2}(u, z)\} = \{\phi_{\ell_1}(u)\varphi_{\ell_2}(z)\}$ are basis functions and $\{\beta_{P,n,i,a,k,c,\ell_1,\ell_2}^f\}, \{\beta_{P,n,j,b,k,c,\ell_1,\ell_2}^g\}$ are coefficients which we can estimate with OLS. The numbers of basis functions for time and the covariate values — denoted by \tilde{c}_n and \tilde{d}_n , respectively — are chosen to increase with the sample size n at some rate. To fix ideas, we will use Legendre polynomials as the basis functions for both the theoretical analysis in this section and the numerical simulations in Section 5. Specifically, for each $\ell_1 \in [\tilde{c}_n]$ and $\ell_2 \in [\tilde{d}_n]$, let the basis functions for time $\{\phi_{\ell_1}(u)\}$ and the covariate values $\{\varphi_{\ell_2}(z)\}$ be *mapped* Legendre polynomials as in Example C.2 and Subsection 3.1.1 of Ding and Zhou [48]. It is straightforward to replace the Legendre polynomials used in our theoretical analysis and simulations with trigonometric polynomials, wavelets, or other Jacobi polynomials.

Next, we introduce the sieve estimators for the time-varying regression functions. Although we do not discuss this topic in detail here, we point interested readers to further discussions of asymptotically optimal linear forecasting for locally stationary processes [45, 86, 35]. We expect similar results for asymptotically optimal nonlinear forecasting to be developed over the next few years.

Remark 4.1 (Sequential sieve estimation). *Our formulation of the sieve estimator from Ding and Zhou [48] accommodates sequential estimation, in the sense that the predictors for rescaled time t/n are only constructed using the information up to rescaled time t/n . We emphasize that sequential estimation is not required for all settings, particularly when certain exogeneity conditions hold. The need for sequential estimation in some settings is due to the martingale difference sequence condition imposed on the error processes in Assumption 3.4 (c.f. Assumption 4.4), which becomes relevant when using our test for variable selection for forecasting and causal inference for time-lagged effects. Note that the same convergence rates are attained whether or not sequential estimation is used, due to the infill asymptotic framework of locally stationary processes. That is, because more observations for each local structure become available as n grows.*

Recall the following notation from Subsection 3.2. Let $\mathfrak{D}_{t,n,i,a}^f, \mathfrak{D}_{t,n,j,b}^g$ be the datasets used to form the estimators $\hat{f}_{t,n,i,a}(t/n, \cdot), \hat{g}_{t,n,j,b}(t/n, \cdot)$ of the time-varying regression functions at rescaled time $t/n \in \mathcal{U}_n$, let $\mathcal{H}_{t,a}^f, \mathcal{H}_{t,b}^g$ be the corresponding input sequences, and let $\mathcal{T}_{t,n,i,a}^f, \mathcal{T}_{t,n,j,b}^g$ be the corresponding sets of times with $T_{t,n,i,a}^f = |\mathcal{T}_{t,n,i,a}^f|, T_{t,n,j,b}^g = |\mathcal{T}_{t,n,j,b}^g|$. Note that each of the estimators $\hat{f}_{t,n,i,a,k,c}(t/n, \cdot), \hat{g}_{t,n,j,b,k,c}(t/n, \cdot)$ of the corresponding time-varying partial response functions at rescaled time $t/n \in \mathcal{U}_n$ may have different numbers of basis functions. Without confusion, we will write the numbers of basis functions as \tilde{c}_n and \tilde{d}_n instead of $\tilde{c}_{t,n,i,a,k,c}^f, \tilde{c}_{t,n,j,b,k,c}^g$ and $\tilde{d}_{t,n,i,a,k,c}^f, \tilde{d}_{t,n,j,b,k,c}^g$ to simplify the presentation below.

For some fixed $n \in \mathbb{N}$, $t \in \mathcal{T}_n$, and $(i, j, a, b) \in \mathcal{D}_n$, denote the design matrices by $\bar{\mathbf{Z}}_{t,n,i,a} \in \mathbb{R}^{T_{t,n,i,a}^f \times \mathbf{d}_Z \tilde{c}_n \tilde{d}_n}$ and $\bar{\mathbf{Z}}_{t,n,j,b} \in \mathbb{R}^{T_{t,n,j,b}^g \times \mathbf{d}_Z \tilde{c}_n \tilde{d}_n}$. The (s, p) -th entries of $\bar{\mathbf{Z}}_{t,n,i,a}$ and $\bar{\mathbf{Z}}_{t,n,j,b}$ are

$$\begin{aligned} \bar{\mathbf{Z}}_{t,n,i,a}^{(s,p)} &= \phi_{\ell_{1,p}}(t_s/n) \varphi_{\ell_{2,p}}(Z_{t_s,n,k_p,c_p}), \\ \bar{\mathbf{Z}}_{t,n,j,b}^{(s,p)} &= \phi_{\ell_{1,p}}(t_s/n) \varphi_{\ell_{2,p}}(Z_{t_s,n,k_p,c_p}), \end{aligned}$$

where we use mappings for the rows $s \mapsto t_s \in \mathcal{T}_{t,n,i,a}^f$ and $s \mapsto t_s \in \mathcal{T}_{t,n,j,b}^g$ which maintain the sequential order of time (i.e. $t_{s_1} < t_{s_2}$ if $s_1 < s_2$), and some mappings for the columns $p \mapsto (k_p, c_p, \ell_{1,p}, \ell_{2,p})$ which determine orderings for the dimension/time-offset/basis-index combinations, where $k_p \in [d_Z]$, $c_p \in C_{k_p}$, $\ell_{1,p} \in [\tilde{c}_n]$, $\ell_{2,p} \in [\tilde{d}_n]$. That is, each row corresponds to one time and each column corresponds to one dimension/time-offset combination with a particular basis-index combination. For

each $n \in \mathbb{N}$, $P \in \mathcal{P}_n$, $i \in [d_X]$, $a \in A_i$, $j \in [d_Y]$, $b \in B_j$ the (time-invariant) coefficient vectors

$$\begin{aligned}\beta_{P,n,i,a}^f &= (\beta_{P,n,i,a,k,c,\ell_1,\ell_2}^f)_{k,c,\ell_1,\ell_2}^\top \in \mathbb{R}^{d_Z \tilde{c}_n \tilde{d}_n}, \\ \beta_{P,n,j,b}^g &= (\beta_{P,n,j,b,k,c,\ell_1,\ell_2}^g)_{k,c,\ell_1,\ell_2}^\top \in \mathbb{R}^{d_Z \tilde{c}_n \tilde{d}_n},\end{aligned}$$

have the following OLS estimators

$$\begin{aligned}\hat{\beta}_{t,n,i,a}^f &= (\bar{\mathbf{Z}}_{t,n,i,a}^\top \bar{\mathbf{Z}}_{t,n,i,a})^{-1} \bar{\mathbf{Z}}_{t,n,i,a}^\top \bar{\mathbf{X}}_{t,n,i,a} = (\hat{\beta}_{t,n,i,a,k,c,\ell_1,\ell_2}^f)_{k,c,\ell_1,\ell_2}^\top \in \mathbb{R}^{d_Z \tilde{c}_n \tilde{d}_n}, \\ \hat{\beta}_{t,n,j,b}^g &= (\bar{\mathbf{Z}}_{t,n,j,b}^\top \bar{\mathbf{Z}}_{t,n,j,b})^{-1} \bar{\mathbf{Z}}_{t,n,j,b}^\top \bar{\mathbf{Y}}_{t,n,j,b} = (\hat{\beta}_{t,n,j,b,k,c,\ell_1,\ell_2}^g)_{k,c,\ell_1,\ell_2}^\top \in \mathbb{R}^{d_Z \tilde{c}_n \tilde{d}_n},\end{aligned}$$

where

$$\bar{\mathbf{X}}_{t,n,i,a} = (X_{t,n,i,a})_{t \in \mathcal{T}_{t,n,i,a}^f}^\top \in \mathbb{R}^{T_{t,n,i,a}^f}, \quad \bar{\mathbf{Y}}_{t,n,j,b} = (Y_{t,n,j,b})_{t \in \mathcal{T}_{t,n,j,b}^g}^\top \in \mathbb{R}^{T_{t,n,j,b}^g}.$$

Finally, the estimators of the time-varying regression functions $f_{P,n,i,a}(t/n, \cdot)$ and $g_{P,n,j,b}(t/n, \cdot)$ at rescaled time $t/n \in \mathcal{U}_n$ are given by

$$\begin{aligned}\hat{f}_{t,n,i,a}(t/n, \cdot) &= \sum_{k=1}^{d_Z} \sum_{c=1}^{C_k} \hat{f}_{t,n,i,a,k,c}(t/n, \cdot), \\ \hat{g}_{t,n,j,b}(t/n, \cdot) &= \sum_{k=1}^{d_Z} \sum_{c=1}^{C_k} \hat{g}_{t,n,j,b,k,c}(t/n, \cdot),\end{aligned}$$

where the estimators of the time-varying partial response functions $f_{P,n,i,a,k,c}(t/n, \cdot)$ and $g_{P,n,j,b,k,c}(t/n, \cdot)$ at rescaled time $t/n \in \mathcal{U}_n$ are given by

$$\begin{aligned}\hat{f}_{t,n,i,a,k,c}(t/n, \cdot) &= \sum_{\ell_1=1}^{\tilde{c}_n} \sum_{\ell_2=1}^{\tilde{d}_n} \hat{\beta}_{t,n,i,a,k,c,\ell_1,\ell_2}^f b_{\ell_1,\ell_2}(t/n, \cdot), \\ \hat{g}_{t,n,j,b,k,c}(t/n, \cdot) &= \sum_{\ell_1=1}^{\tilde{c}_n} \sum_{\ell_2=1}^{\tilde{d}_n} \hat{\beta}_{t,n,j,b,k,c,\ell_1,\ell_2}^g b_{\ell_1,\ell_2}(t/n, \cdot).\end{aligned}$$

Although we only discuss the sieve estimator here, we emphasize that any black-box time-varying regression estimator can be used with the dGCM test. For example, we can use time-varying regression estimators based on kernel smoothing [161, 186, 177, 27, 48]. To use kernel smoothing estimators for sequential estimation, we can use one-sided temporal kernels so that observations after rescaled time t/n receive a weight of zero. This is practically important because “local” nonparametric estimators are naturally far more computationally efficient for sequential estimation than “global” nonparametric estimators in the absence of efficient online estimation procedures for the latter.

4.4 Locally stationary error processes

We will now introduce the causal representations of the locally stationary error processes from Subsection 4.3. For each $a \in A$, $b \in B$, define the input sequences

$$\mathcal{H}_{t,a}^\varepsilon = (\eta_{t,a}^\varepsilon, \eta_{t,a-1}^\varepsilon, \dots), \quad \mathcal{H}_{t,b}^\xi = (\eta_{t,b}^\xi, \eta_{t,b-1}^\xi, \dots),$$

where $(\eta_{t,a}^\varepsilon, \eta_{t,b}^\xi)_{t \in \mathbb{Z}}$ is a sequence of iid random vectors. Denote the dimension of $\eta_{t,a}^\varepsilon = \eta_{t,a,n}^\varepsilon$ by $d_\varepsilon^\eta = d_{\varepsilon,n}^\eta$ and the dimension of $\eta_{t,b}^\xi = \eta_{t,b,n}^\xi$ by $d_\xi^\eta = d_{\xi,n}^\eta$, both of which can change with n . For the next assumption, let $\mathcal{H}_t^f = (\mathcal{H}_{t,a}^f)_{a \in A}$, $\mathcal{H}_t^g = (\mathcal{H}_{t,b}^g)_{b \in B}$ and $\mathcal{H}_{t,a}^f = (\mathcal{H}_{t,a}^{\mathfrak{D}^f}, \mathcal{H}_t^Z)$, $\mathcal{H}_{t,b}^g = (\mathcal{H}_{t,b}^{\mathfrak{D}^g}, \mathcal{H}_t^Z)$, where the input sequences $\mathcal{H}_{t,a}^{\mathfrak{D}^f}$, $\mathcal{H}_{t,b}^{\mathfrak{D}^g}$ were defined in Subsection 4.3 and \mathcal{H}_t^Z was defined in Subsection 4.2.

Assumption 4.4 (Causal representations of the error processes). *Assume that for each $n \in \mathbb{N}$, $P \in \mathcal{P}_n$, $(i, j, a, b) \in \mathcal{D}_n$, $t \in \mathcal{T}_n$, the error processes from Subsection 4.3 can be represented as*

$$\varepsilon_{P,t,n,i,a} = \tilde{G}_{P,n,i,a}^\varepsilon(t/n, \mathcal{H}_{t,a}^\varepsilon), \quad \xi_{P,t,n,j,b} = \tilde{G}_{P,n,j,b}^\xi(t/n, \mathcal{H}_{t,b}^\xi),$$

with $\mathbb{E}_P(\varepsilon_{P,t,n,i,a}|\mathcal{H}_t^{\hat{g}}) = 0$ and $\mathbb{E}_P(\xi_{P,t,n,j,b}|\mathcal{H}_t^{\hat{f}}) = 0$, where the input sequences $\mathcal{H}_t^{\hat{g}}$, $\mathcal{H}_t^{\hat{f}}$ were defined above. The causal representations

$$\tilde{\varepsilon}_{P,t,n,i,a}(u) = \tilde{G}_{P,n,i,a}^{\varepsilon}(u, \mathcal{H}_{t,a}^{\varepsilon}), \quad \tilde{\xi}_{P,t,n,j,b}(u) = \tilde{G}_{P,n,j,b}^{\xi}(u, \mathcal{H}_{t,b}^{\xi}),$$

are defined at all $u \in \mathcal{U}_n$, so that we have $\varepsilon_{P,t,n,i,a} = \tilde{\varepsilon}_{P,t,n,i,a}(t/n)$, $\xi_{P,t,n,j,b} = \tilde{\xi}_{P,t,n,j,b}(t/n)$. $\tilde{G}_{P,n,i,a}^{\varepsilon}(u, \cdot)$ and $\tilde{G}_{P,n,j,b}^{\xi}(u, \cdot)$ are measurable functions from $(\mathbb{R}^{d_{\varepsilon}})^{\infty}$ and $(\mathbb{R}^{d_{\xi}})^{\infty}$, respectively, to \mathbb{R} — where we endow $(\mathbb{R}^{d_{\varepsilon}})^{\infty}$ and $(\mathbb{R}^{d_{\xi}})^{\infty}$ with the σ -algebra generated by all finite projections — so that $\tilde{G}_{P,n,i,a}^{\varepsilon}(u, \mathcal{H}_{s,a}^{\varepsilon})$, $\tilde{G}_{P,n,j,b}^{\xi}(u, \mathcal{H}_{s,b}^{\xi})$ are well-defined random variables for each $s \in \mathbb{Z}$ and $(\tilde{G}_{P,n,i,a}^{\varepsilon}(u, \mathcal{H}_{s,a}^{\varepsilon}))_{s \in \mathbb{Z}}$, $(\tilde{G}_{P,n,j,b}^{\xi}(u, \mathcal{H}_{s,b}^{\xi}))_{s \in \mathbb{Z}}$ are stationary ergodic processes.

As in Subsection 3.3, we have not defined the input sequences for the error processes separately for each dimension, because without loss of generality we may define the measurable functions $\tilde{G}_{P,n,i,a}^{\varepsilon}(u, \cdot)$, $\tilde{G}_{P,n,j,b}^{\xi}(u, \cdot)$ and inputs $\eta_{t,a}^{\varepsilon}, \eta_{t,b}^{\xi}$ so that each dimension of the error processes has idiosyncratic inputs.

Using the causal representations of the univariate error processes, we have the following causal representations of the vector-valued error processes

$$\begin{aligned} \tilde{\varepsilon}_{P,t,n}(u) &= \tilde{\mathbf{G}}_{P,n}^{\varepsilon}(u, \mathcal{H}_t^{\varepsilon}) = (\tilde{G}_{P,n,i,a}^{\varepsilon}(u, \mathcal{H}_{t,a}^{\varepsilon}))_{i \in [d_X], a \in A_i}, \\ \tilde{\xi}_{P,t,n}(u) &= \tilde{\mathbf{G}}_{P,n}^{\xi}(u, \mathcal{H}_t^{\xi}) = (\tilde{G}_{P,n,j,b}^{\xi}(u, \mathcal{H}_{t,b}^{\xi}))_{j \in [d_Y], b \in B_j}, \end{aligned}$$

so that we have $\varepsilon_{P,t,n} = \tilde{\varepsilon}_{P,t,n}(t/n)$, $\xi_{P,t,n} = \tilde{\xi}_{P,t,n}(t/n)$, where $\mathcal{H}_t^{\varepsilon} = (\eta_t^{\varepsilon}, \eta_{t-1}^{\varepsilon}, \dots)$, $\mathcal{H}_t^{\xi} = (\eta_t^{\xi}, \eta_{t-1}^{\xi}, \dots)$ with $\eta_t^{\varepsilon} = (\eta_{t,a}^{\varepsilon})_{a \in A}$, $\eta_t^{\xi} = (\eta_{t,b}^{\xi})_{b \in B}$ for each $t \in \mathbb{Z}$. Similarly, for each $m = (i, j, a, b) \in \mathcal{D}_n$ the error products can be represented as

$$\tilde{R}_{P,t,n,m}(u) = \tilde{G}_{P,n,m}^R(u, \mathcal{H}_{t,m}^R) = \tilde{G}_{P,n,i,a}^{\varepsilon}(u, \mathcal{H}_{t,a}^{\varepsilon}) \tilde{G}_{P,n,j,b}^{\xi}(u, \mathcal{H}_{t,b}^{\xi}),$$

so that we have $R_{P,t,n,m} = \tilde{R}_{P,t,n,m}(t/n)$, where $\mathcal{H}_{t,m}^R = (\eta_{t,m}^R, \eta_{t-1,m}^R, \dots)$ with $\eta_{t,m}^R = (\eta_{t,a}^{\varepsilon}, \eta_{t,b}^{\xi})^{\top}$ for each $t \in \mathbb{Z}$. Also, we have the following causal representation of the nonstationary \mathbb{R}^{D_n} -valued process of all the products of errors

$$\tilde{\mathbf{R}}_{P,n,t}(u) = \tilde{\mathbf{G}}_{P,n}^R(u, \mathcal{H}_t^R) = (\tilde{G}_{P,n,m}^R(u, \mathcal{H}_{t,m}^R))_{m \in \mathcal{D}_n},$$

so that we have $\mathbf{R}_{P,t,n} = \tilde{\mathbf{R}}_{P,n,t}(t/n)$, where $\mathcal{H}_t^R = (\eta_t^R, \eta_{t-1}^R, \dots)$ and $\eta_t^R = (\eta_t^{\varepsilon}, \eta_t^{\xi})^{\top}$ for each $t \in \mathbb{Z}$. We emphasize that for a fixed $P \in \mathcal{P}_n$, $u \in \mathcal{U}_n$, and $n \in \mathbb{N}$, we have that $\tilde{\mathbf{G}}_{P,n}^R(u, \mathcal{H}_s^R)$ is a well-defined random vector for each $s \in \mathbb{Z}$ and $(\tilde{\mathbf{G}}_{P,n}^R(u, \mathcal{H}_s^R))_{s \in \mathbb{Z}}$ is a stationary ergodic \mathbb{R}^{D_n} -valued process.

4.5 Assumptions on dependence and nonstationarity

In this subsection, we impose assumptions on the rate of decay in temporal dependence and the degree of nonstationarity of the observed processes and error processes. We emphasize that the assumptions here are strictly stronger than those in Subsection 3.4. We impose these stronger assumptions to guarantee that the sieve time-varying nonlinear regression estimator achieves the convergence rates required by Theorem 3.1. Note that the assumptions here require that the nonstationary processes evolve “smoothly” in time, which excludes nonstationary processes with abrupt changes. We do this mainly to simplify the presentation, and we discuss extensions to nonstationary processes with *both* smooth and abrupt changes in Subsection A.6.

Denote the set of well-defined tuples of observed processes, dimensions, and time-offsets by

$$\mathbb{W} = \{(X, i, a) : i \in [d_X], a \in A_i\} \cup \{(Y, j, b) : j \in [d_Y], b \in B_j\} \cup \{(Z, k, c) : k \in [d_Z], c \in C_k\},$$

so that we may conveniently refer to such well-defined combinations by $(W, l, d) \in \mathbb{W}$. Also, denote the set of well-defined tuples of error processes, dimensions, and time-offsets by

$$\mathbb{E} = \{(\varepsilon, i, a) : i \in [d_X], a \in A_i\} \cup \{(\xi, j, b) : j \in [d_Y], b \in B_j\},$$

so that we may write $(e, l, d) \in \mathbb{E}$ to refer to any such combination.

Again, we quantify temporal dependence via the functional dependence measure of Wu [170]. Let $(\tilde{\eta}_t^X, \tilde{\eta}_t^Y, \tilde{\eta}_t^Z)_{t \in \mathbb{Z}}$ be an iid copy of $(\eta_t^X, \eta_t^Y, \eta_t^Z)_{t \in \mathbb{Z}}$. Going forward, the inputs with the tilde are from $(\tilde{\eta}_t^X, \tilde{\eta}_t^Y, \tilde{\eta}_t^Z)_{t \in \mathbb{Z}}$. For any tuple $(W, l, d) \in \mathbb{W}$ corresponding to a well-defined combination of an observed process, dimension, and time-offset, define

$$\tilde{\mathcal{H}}_{t,d,h}^W = (\eta_{t+d}^W, \dots, \eta_{t-h+1+d}^W, \tilde{\eta}_{t-h+d}^W, \eta_{t-h-1+d}^W, \dots)$$

to be $\mathcal{H}_{t,d}^W$ with the input η_{t-h+d}^W replaced with the iid copy $\tilde{\eta}_{t-h+d}^W$. For example, for $i \in [d_X]$, $a \in A_i$, we have that $\tilde{\eta}_{t-h+a}^X$ is the copy of the input η_{t-h+a}^X in the input sequence $\mathcal{H}_{t,a}^X$ used in the causal representation of $X_{t,n,i,a}$. Analogously, for $\mathbf{W} \in \{\mathbf{X}, \mathbf{Y}, \mathbf{Z}\}$ define $\tilde{\mathcal{H}}_{t,h}^{\mathbf{W}}$ as $\mathcal{H}_t^{\mathbf{W}}$ with the input $\eta_{t-h}^{\mathbf{W}}$ replaced with the iid copy $\tilde{\eta}_{t-h}^{\mathbf{W}}$ as in Subsection 4.2.

For any tuple $(e, l, d) \in \mathbb{E}$ corresponding to a well-defined combination of an error process, dimension, and time-offset, define

$$\tilde{\mathcal{H}}_{t,d,h}^e = (\eta_{t,d}^e, \dots, \eta_{t-h+1,d}^e, \tilde{\eta}_{t-h,d}^e, \eta_{t-h-1,d}^e, \dots)$$

to be $\mathcal{H}_{t,d}^e$ with the input $\eta_{t-h,d}^e$ replaced with the iid copy $\tilde{\eta}_{t-h,d}^e$. Analogously, for $e \in \{\varepsilon, \xi\}$ define $\tilde{\mathcal{H}}_{t,h}^e$ as \mathcal{H}_t^e with the input η_{t-h}^e replaced with the iid copy $\tilde{\eta}_{t-h}^e$ as in Subsection 4.4. Also, for the product of errors define $\tilde{\mathcal{H}}_{t,m,h}^R$ as $\mathcal{H}_{t,m}^R$ with the input $\eta_{t-h,m}^R$ replaced with the iid copy $\tilde{\eta}_{t-h,m}^R$ for $m = (i, j, a, b) \in \mathcal{D}_n$. Analogously, define $\tilde{\mathcal{H}}_{t,h}^R$ as \mathcal{H}_t^R with the input η_{t-h}^R replaced with the iid copy $\tilde{\eta}_{t-h}^R$ as in Subsection 4.4. Now, we define the functional dependence measures of the processes.

Definition 4.1 (Functional dependence measures). *We define the following measures of temporal dependence for each $n \in \mathbb{N}$, $P \in \mathcal{P}_n$, $u \in \mathcal{U}_n$, $t \in \mathcal{T}_n$. First, define the functional dependence measures of the observed processes $\tilde{G}_{n,l,d}^W(u, \mathcal{H}_{t,d}^W)$ for each $(W, l, d) \in \mathbb{W}$ with $h \in \mathbb{N}_0$, and some $q \geq 1$ as*

$$\theta_{P,u,t,n,l,d}^W(h, q) = [\mathbb{E}_P(|\tilde{G}_{n,l,d}^W(u, \mathcal{H}_{t,d}^W) - \tilde{G}_{n,l,d}^W(u, \tilde{\mathcal{H}}_{t,d,h}^W)|^q)]^{1/q},$$

and for the vector-valued process $\tilde{\mathbf{G}}_n^{\mathbf{W}}(u, \mathcal{H}_t^{\mathbf{W}})$ for each $\mathbf{W} \in \{\mathbf{X}, \mathbf{Y}, \mathbf{Z}\}$ with $h \in \mathbb{N}_0$, and some $q \geq 1$, $r \geq 1$ as

$$\theta_{P,u,t,n}^{\mathbf{W}}(h, q, r) = [\mathbb{E}_P(\|\tilde{\mathbf{G}}_n^{\mathbf{W}}(u, \mathcal{H}_t^{\mathbf{W}}) - \tilde{\mathbf{G}}_n^{\mathbf{W}}(u, \tilde{\mathcal{H}}_{t,h}^{\mathbf{W}}\|_r^q)]^{1/q}.$$

Second, define the L^∞ versions of the functional dependence measures of the error processes $\tilde{G}_{P,n,l,d}^e(u, \mathcal{H}_{t,d}^e)$ for each $(e, l, d) \in \mathbb{E}$ with $h \in \mathbb{N}_0$ as

$$\theta_{P,u,t,n,l,d}^{e,\infty}(h) = \inf\{K \geq 0 : \mathbb{P}_P(|\tilde{G}_{P,n,l,d}^e(u, \mathcal{H}_{t,d}^e) - \tilde{G}_{P,n,l,d}^e(u, \tilde{\mathcal{H}}_{t,d,h}^e)| > K) = 0\},$$

and for the vector-valued process $\tilde{\mathbf{G}}_{P,n}^e(u, \mathcal{H}_t^e)$ for each $e \in \{\varepsilon, \xi\}$ with $h \in \mathbb{N}_0$, and some $r \geq 1$ as

$$\theta_{P,u,t,n}^{e,\infty}(h, r) = \inf\{K \geq 0 : \mathbb{P}_P(\|\tilde{\mathbf{G}}_{P,n}^e(u, \mathcal{H}_t^e) - \tilde{\mathbf{G}}_{P,n}^e(u, \tilde{\mathcal{H}}_{t,h}^e)\|_r > K) = 0\}.$$

Third, define the functional dependence measures of the processes of error products $\tilde{G}_{P,n,m}^R(u, \mathcal{H}_{t,m}^R)$ for each $m = (i, j, a, b) \in \mathcal{D}_n$ with $h \in \mathbb{N}_0$, and some $q \geq 1$ as

$$\theta_{P,u,t,n,m}^R(h, q) = [\mathbb{E}_P(|\tilde{G}_{P,n,m}^R(u, \mathcal{H}_{t,m}^R) - \tilde{G}_{P,n,m}^R(u, \tilde{\mathcal{H}}_{t,m,h}^R)|^q)]^{1/q},$$

and for the vector-valued process $\tilde{\mathbf{G}}_{P,n}^R(u, \mathcal{H}_t^R)$ with $h \in \mathbb{N}_0$, and some $q \geq 1$, $r \geq 1$ as

$$\theta_{P,u,t,n}^R(h, q, r) = [\mathbb{E}_P(\|\tilde{\mathbf{G}}_{P,n}^R(u, \mathcal{H}_t^R) - \tilde{\mathbf{G}}_{P,n}^R(u, \tilde{\mathcal{H}}_{t,h}^R)\|_r^q)]^{1/q}.$$

Next, we introduce an assumption imposing a uniform polynomial decay of the temporal dependence. Note that we will often write the time as 0 when the time of the input sequence does not matter because of stationarity. For some sequence of collections of distributions $(\mathcal{P}_n)_{n \in \mathbb{N}}$, we make the following assumption.

Assumption 4.5 (Distribution-uniform decay of temporal dependence). *Assume that there exist $\bar{\Theta} > 0$, $\bar{\beta} > 2$, $\bar{q} > 4$, such that for all $n \in \mathbb{N}$, $u \in \mathcal{U}_n$, and observed processes $(W, l, d) \in \mathbb{W}$, it holds that*

$$\sup_{P \in \mathcal{P}_n} [\mathbb{E}_P(|\tilde{G}_{n,l,d}^W(u, \mathcal{H}_{0,d}^W)|^{\bar{q}})]^{1/\bar{q}} \leq \bar{\Theta}, \quad \sup_{P \in \mathcal{P}_n} \theta_{P,u,0,n,l,d}^W(h, \bar{q}) \leq \bar{\Theta} \cdot (h \vee 1)^{-\bar{\beta}}, \quad h \geq 0.$$

Also, assume that there exist $\bar{\Theta}^\infty > 0$, $\bar{\beta}^\infty > 2$, such that for all $n \in \mathbb{N}$, $u \in \mathcal{U}_n$, and error processes $(e, l, d) \in \mathbb{E}$, it holds that

$$\sup_{P \in \mathcal{P}_n} \|G_{P,n,l,d}^e(u, \mathcal{H}_{0,d}^e)\|_{L^\infty(P)} \leq \bar{\Theta}^\infty, \quad \sup_{P \in \mathcal{P}_n} \theta_{P,u,0,n,l,d}^{e,\infty}(h) \leq \bar{\Theta}^\infty \cdot (h \vee 1)^{-\bar{\beta}^\infty}, \quad h \geq 0.$$

For additional control in terms of the product of errors alone, also assume that there exist $\bar{\Theta}^R > 0$, $\bar{\beta}^R > 3$, $\bar{q}^R > 4$, such that for all $n \in \mathbb{N}$, $u \in \mathcal{U}_n$, $m = (i, j, a, b) \in \mathcal{D}_n$, it holds that

$$\sup_{P \in \mathcal{P}_n} [\mathbb{E}_P(|\tilde{G}_{P,n,m}^R(u, \mathcal{H}_{0,m}^R)|^{\bar{q}^R})]^{1/\bar{q}^R} \leq \bar{\Theta}^R, \quad \sup_{P \in \mathcal{P}_n} \theta_{P,u,0,n,m}^R(h, \bar{q}^R) \leq \bar{\Theta}^R \cdot (h \vee 1)^{-\bar{\beta}^R}, \quad h \geq 0.$$

In view of Assumption 4.5, we have the following bounds on the functional dependence measures of the corresponding vector-valued processes for each $n \in \mathbb{N}$, $u \in \mathcal{U}_n$ by Jensen's inequality. For the vector-valued process of error products, we have

$$\sup_{P \in \mathcal{P}_n} [\mathbb{E}_P(|\tilde{G}_{P,n}^R(u, \mathcal{H}_0^R)|_2^{\bar{q}^R})]^{1/\bar{q}^R} \leq D_n^{\frac{1}{2}} \bar{\Theta}^R, \quad \sup_{P \in \mathcal{P}_n} \theta_{P,u,0,n}^R(h, \bar{q}^R, 2) \leq D_n^{\frac{1}{2}} \bar{\Theta}^R \cdot (h \vee 1)^{-\bar{\beta}^R}, \quad h \geq 0.$$

Also, for each of the vector-valued observed processes $\mathbf{W} \in (\mathbf{X}, \mathbf{Y}, \mathbf{Z})$, we have

$$\sup_{P \in \mathcal{P}_n} [\mathbb{E}_P(|\tilde{\mathbf{G}}_n^{\mathbf{W}}(u, \mathcal{H}_0^{\mathbf{W}})|_2^{\bar{q}})]^{1/\bar{q}} \leq D_n^{\frac{1}{2}} \bar{\Theta}, \quad \sup_{P \in \mathcal{P}_n} \theta_{P,u,0,n}^{\mathbf{W}}(h, \bar{q}, 2) \leq D_n^{\frac{1}{2}} \bar{\Theta} \cdot (h \vee 1)^{-\bar{\beta}}, \quad h \geq 0.$$

Lastly, for each of the vector-valued error processes $\mathbf{e} \in (\varepsilon, \xi)$, we have

$$\sup_{P \in \mathcal{P}_n} \left\| |\tilde{\mathbf{G}}_{P,n}^{\mathbf{e}}(u, \mathcal{H}_0^{\mathbf{e}})|_2 \right\|_{L^\infty(P)} \leq D_n^{\frac{1}{2}} \bar{\Theta}^\infty, \quad \sup_{P \in \mathcal{P}_n} \theta_{P,u,0,n}^{\mathbf{e},\infty}(h, 2) \leq D_n^{\frac{1}{2}} \bar{\Theta}^\infty \cdot (h \vee 1)^{-\bar{\beta}^\infty}, \quad h \geq 0.$$

Next, we discuss an additional regularity condition required by the sieve estimator that is analogous to Lemma 3.1 in Ding and Zhou [48]. Recall the set of basis functions $\{\varphi_{\ell_2}(z)\}$ from Subsection 4.3. For each $n \in \mathbb{N}$, $P \in \mathcal{P}_n$, $t \in \mathcal{T}_n$, $i \in [d_X]$, $a \in A_i$, $j \in [d_Y]$, $b \in B_j$ let

$$\begin{aligned} \mathbf{w}_{t,n}^{\varphi(Z)} &= (\varphi_{\ell_2}(Z_{t,n,k,c}))_{k \in [d_Z], c \in C_k, 1 \leq \ell_2 \leq \tilde{d}_n}, \\ \mathbf{w}_{P,t,n,i,a}^{\varphi(Z),\varepsilon} &= (\varphi_{\ell_2}(Z_{t,n,k,c}) \varepsilon_{P,t,n,i,a})_{k \in [d_Z], c \in C_k, 1 \leq \ell_2 \leq \tilde{d}_n}, \\ \mathbf{w}_{P,t,n,j,b}^{\varphi(Z),\xi} &= (\varphi_{\ell_2}(Z_{t,n,k,c}) \xi_{P,t,n,j,b})_{k \in [d_Z], c \in C_k, 1 \leq \ell_2 \leq \tilde{d}_n}. \end{aligned}$$

As in Subsection 3.2 in Ding and Zhou [48], the $\mathbb{R}^{d_Z \tilde{d}_n}$ -valued processes $\mathbf{w}_{t,n}^{\varphi(Z)}$, $\mathbf{w}_{P,t,n,i,a}^{\varphi(Z),\varepsilon}$, and $\mathbf{w}_{P,t,n,j,b}^{\varphi(Z),\xi}$ all have causal representations

$$\begin{aligned} \mathbf{w}_{t,n}^{\varphi(Z)} &= \tilde{\mathbf{G}}_n^{\varphi(Z)}(t/n, \mathcal{H}_t^{\varphi(Z)}), \\ \mathbf{w}_{P,t,n,i,a}^{\varphi(Z),\varepsilon} &= \tilde{\mathbf{G}}_{P,n,i,a}^{\varphi(Z),\varepsilon}(t/n, \mathcal{H}_{t,a}^{\varphi(Z),\varepsilon}), \\ \mathbf{w}_{P,t,n,j,b}^{\varphi(Z),\xi} &= \tilde{\mathbf{G}}_{P,n,j,b}^{\varphi(Z),\xi}(t/n, \mathcal{H}_{t,b}^{\varphi(Z),\xi}), \end{aligned}$$

where

$$\begin{aligned} \mathcal{H}_t^{\varphi(Z)} &= (\eta_t^{\varphi(Z)}, \eta_{t-1}^{\varphi(Z)}, \dots), \\ \mathcal{H}_{t,a}^{\varphi(Z),\varepsilon} &= (\eta_{t,a}^{\varphi(Z),\varepsilon}, \eta_{t-1,a}^{\varphi(Z),\varepsilon}, \dots), \\ \mathcal{H}_{t,b}^{\varphi(Z),\xi} &= (\eta_{t,b}^{\varphi(Z),\xi}, \eta_{t-1,b}^{\varphi(Z),\xi}, \dots), \end{aligned}$$

with $\eta_t^{\varphi(Z)} = \eta_{t+c_{\max}}^Z$, $\eta_{t,a}^{\varphi(Z),\varepsilon} = (\eta_{t+c_{\max}}^Z, \eta_{t+a}^X)^\top$, and $\eta_{t,b}^{\varphi(Z),\xi} = (\eta_{t+c_{\max}}^Z, \eta_{t+b}^Y)^\top$.

Define the functional dependence measures of the vector-valued processes $\mathbf{w}_{t,n}^{\varphi(Z)}$, $\mathbf{w}_{P,t,n,i,a}^{\varphi(Z),\varepsilon}$, $\mathbf{w}_{P,t,n,j,b}^{\varphi(Z),\xi}$ by

$$\begin{aligned} \theta_{P,u,t,n}^{\varphi(Z)}(h, q, 2) &= [\mathbb{E}_P(|\tilde{\mathbf{G}}_n^{\varphi(Z)}(u, \mathcal{H}_t^{\varphi(Z)}) - \tilde{\mathbf{G}}_n^{\varphi(Z)}(u, \tilde{\mathcal{H}}_{t,h}^{\varphi(Z)})|_2^q)]^{1/q}, \\ \theta_{P,u,t,n,i,a}^{\varphi(Z),\varepsilon}(h, q, 2) &= [\mathbb{E}_P(|\tilde{\mathbf{G}}_{P,n,i,a}^{\varphi(Z),\varepsilon}(u, \mathcal{H}_{t,a}^{\varphi(Z),\varepsilon}) - \tilde{\mathbf{G}}_{P,n,i,a}^{\varphi(Z),\varepsilon}(u, \tilde{\mathcal{H}}_{t,a,h}^{\varphi(Z),\varepsilon})|_2^q)]^{1/q}, \\ \theta_{P,u,t,n,j,b}^{\varphi(Z),\xi}(h, q, 2) &= [\mathbb{E}_P(|\tilde{\mathbf{G}}_{P,n,j,b}^{\varphi(Z),\xi}(u, \mathcal{H}_{t,b}^{\varphi(Z),\xi}) - \tilde{\mathbf{G}}_{P,n,j,b}^{\varphi(Z),\xi}(u, \tilde{\mathcal{H}}_{t,b,h}^{\varphi(Z),\xi})|_2^q)]^{1/q}. \end{aligned}$$

Recall $\bar{\Theta}$, $\bar{\Theta}^\infty$, $\bar{\beta}$, $\bar{\beta}^\infty$, and \bar{q} from Assumption 4.5. Using the same arguments from Lemma 3.1 from Ding and Zhou [48], for all $n \in \mathbb{N}$, $u \in \mathcal{U}_n$, the vector-valued processes $\mathbf{w}_{P,t,n,i,a}^{\varphi(Z),\varepsilon}$, $\mathbf{w}_{P,t,n,j,b}^{\varphi(Z),\xi}$ satisfy

$$\sup_{P \in \mathcal{P}_n} [\mathbb{E}_P(\|\tilde{\mathbf{G}}_{P,n,i,a}^{\varphi(Z),\varepsilon}(u, \mathcal{H}_{t,a}^{\varphi(Z),\varepsilon})\|_2^{\bar{q}})]^{1/\bar{q}} \leq D_n^{\frac{1}{2}} \tilde{\Theta}, \quad \sup_{P \in \mathcal{P}_n} \theta_{P,u,t,n,i,a}^{\varphi(Z),\varepsilon}(h, \bar{q}, 2) \leq D_n^{\frac{1}{2}} \tilde{\Theta} \cdot (h \vee 1)^{-\bar{\beta}},$$

$$\sup_{P \in \mathcal{P}_n} [\mathbb{E}_P(\|\tilde{\mathbf{G}}_{P,n,j,b}^{\varphi(Z),\xi}(u, \mathcal{H}_{t,b}^{\varphi(Z),\xi})\|_2^{\bar{q}})]^{1/\bar{q}} \leq D_n^{\frac{1}{2}} \tilde{\Theta}, \quad \sup_{P \in \mathcal{P}_n} \theta_{P,u,t,n,j,b}^{\varphi(Z),\xi}(h, \bar{q}, 2) \leq D_n^{\frac{1}{2}} \tilde{\Theta} \cdot (h \vee 1)^{-\bar{\beta}},$$

for $h \geq 0$, where $\bar{q} = \bar{q} > 4$ with $\tilde{\beta} = \min(\bar{\beta}, \bar{\beta}^\infty) > 2$ and $\tilde{\Theta} = 2K_1(\max(\bar{\Theta}^\infty, \bar{\Theta}))^2 > 0$ where the constant factor $K_1 > 0$ is due to the basis functions. Similarly, for all $n \in \mathbb{N}$, $u \in \mathcal{U}_n$, the vector-valued process $\mathbf{w}_{t,n}^{\varphi(Z)}$ satisfies

$$\sup_{P \in \mathcal{P}_n} [\mathbb{E}_P(\|\tilde{\mathbf{G}}_n^{\varphi(Z)}(u, \mathcal{H}_t^{\varphi(Z)})\|_2^{\bar{q}})]^{1/\bar{q}} \leq D_n^{\frac{1}{2}} \bar{\Theta} K_2, \quad \sup_{P \in \mathcal{P}_n} \theta_{P,u,t,n}^{\varphi(Z)}(h, \bar{q}, 2) \leq D_n^{\frac{1}{2}} \bar{\Theta} K_2 \cdot (h \vee 1)^{-\bar{\beta}},$$

for $h \geq 0$, where the constant factor $K_2 > 0$ is due to the basis functions.

For Theorem 3.1, we only require that the total variation of the causal mechanism of the process of error products can be bounded distribution-uniformly. However, the sieve estimator requires the stronger assumption that the causal mechanisms of the observed processes and error processes are stochastic Lipschitz functions of rescaled time. We impose the following regularity conditions to control the nonstationarity uniformly over a sequence of collections of distributions $(\mathcal{P}_n)_{n \in \mathbb{N}}$.

Assumption 4.6 (Distribution-uniform stochastic Lipschitz condition for nonstationarity). *For each $n \in \mathbb{N}$, $(W, l, d) \in \mathbb{W}$, $(e, l, d) \in \mathbb{E}$, and $t \in \mathbb{Z}$, we assume that $\tilde{G}_{n,l,d}^W(\cdot, \mathcal{H}_{t,d}^W)$ and $\tilde{G}_{P,n,l,d}^e(\cdot, \mathcal{H}_{t,d}^e)$ are stochastic Lipschitz functions of rescaled time $u \in \mathcal{U}_n$. Recall $\bar{\Theta} > 0$, $\bar{q} > 4$ from Assumption 4.5. Assume that there exists a constant $\bar{L} > 0$ such that for all $n \in \mathbb{N}$, $u, v \in \mathcal{U}_n$, $(W, l, d) \in \mathbb{W}$, $(e, l, d) \in \mathbb{E}$, it holds that*

$$\sup_{P \in \mathcal{P}_n} [\mathbb{E}_P(|\tilde{G}_{n,l,d}^W(u, \mathcal{H}_{0,d}^W) - \tilde{G}_{n,l,d}^W(v, \mathcal{H}_{0,d}^W)|^{\bar{q}})]^{1/\bar{q}} \leq \bar{L} \bar{\Theta} |u - v|,$$

$$\sup_{P \in \mathcal{P}_n} [\mathbb{E}_P(|\tilde{G}_{P,n,l,d}^e(u, \mathcal{H}_{0,d}^e) - \tilde{G}_{P,n,l,d}^e(v, \mathcal{H}_{0,d}^e)|^{\bar{q}})]^{1/\bar{q}} \leq \bar{L} \bar{\Theta} |u - v|.$$

In view of Assumption 4.6, there exist $\tilde{L}^R = \bar{L} > 0$, $\tilde{q}^R = \bar{q} > 4$, $\tilde{\Theta}^R = 2(\max(\bar{\Theta}^\infty, \bar{\Theta}))^2$ such that for all $n \in \mathbb{N}$, $u, v \in \mathcal{U}_n$, $m = (i, j, a, b) \in \mathcal{D}_n$ we have

$$\sup_{P \in \mathcal{P}_n} [\mathbb{E}_P(|\tilde{G}_{P,n,m}^R(u, \mathcal{H}_{0,m}^R) - \tilde{G}_{P,n,m}^R(v, \mathcal{H}_{0,m}^R)|^{\tilde{q}^R})]^{1/\tilde{q}^R} \leq \tilde{L}^R \tilde{\Theta}^R |u - v|.$$

This follows from adding and subtracting cross-terms, the triangle inequality, the distributive property, Hölder's inequality, and applying the moment bounds and stochastic Lipschitz conditions for the individual error processes from Assumptions 4.5 and 4.6. It is easy to verify that Assumption 3.6 is satisfied under this stronger condition on the nonstationarity.

Also, using the same arguments as Lemma 3.1 from Ding and Zhou [48], the individual dimensions of the vector-valued processes $\mathbf{w}_{P,t,n,i,a}^{\varphi(Z),\varepsilon}$ and $\mathbf{w}_{P,t,n,j,b}^{\varphi(Z),\xi}$ can be shown to satisfy this stochastic Lipschitz condition for moment $\bar{q}/2 > 2$ with Lipschitz constant $2K_1 \bar{L}(\max(\bar{\Theta}^\infty, \bar{\Theta}))^2 > 0$, where the constant factor $K_1 > 0$ is due to the basis functions. Similarly, the individual dimensions of the vector-valued process $\mathbf{w}_{t,n}^{\varphi(Z)}$ can be shown to satisfy this stochastic Lipschitz condition for moment $\bar{q} > 4$ with Lipschitz constant $K_2 \bar{L} \bar{\Theta} > 0$, where the constant factor $K_2 > 0$ is due to the basis functions.

4.6 Assumptions on local long-run covariances

To ensure fast convergence rates by the sieve estimator, we require the following assumptions on the local long-run covariance matrices. Note that these assumptions are not made in Section 3.

Definition 4.2 (Local long-run covariance matrices of error products). *For each $n \in \mathbb{N}$, $P \in \mathcal{P}_n$, $u \in \mathcal{U}_n$, define the local long-run covariance matrix $\tilde{\Sigma}_{P,n}^R(u) \in \mathbb{R}^{D_n \times D_n}$ for the \mathbb{R}^{D_n} -valued stationary process $(\tilde{\mathbf{G}}_{P,n}^R(u, \mathcal{H}_t^R))_{t \in \mathbb{Z}}$ by*

$$\tilde{\Sigma}_{P,n}^R(u) = \sum_{h \in \mathbb{Z}} \text{Cov}_P(\tilde{\mathbf{G}}_{P,n}^R(u, \mathcal{H}_0^R), \tilde{\mathbf{G}}_{P,n}^R(u, \mathcal{H}_h^R)).$$

By Lemma B.5, the local long-run covariance matrices of $\mathbf{w}_{t,n}^{\varphi(Z)}$, $\mathbf{w}_{P,t,n,i,a}^{\varphi(Z),\varepsilon}$, $\mathbf{w}_{P,t,n,j,b}^{\varphi(Z),\xi}$ are well-defined in view of the discussion following Assumption 4.5. Now, we will define the local long-run and integrated long-run covariance matrices of these processes as in Subsection 3.2 of Ding and Zhou [48].

Definition 4.3 (Local long-run and integrated long-run covariance matrices). *For each $n \in \mathbb{N}$, $P \in \mathcal{P}_n$, $u \in \mathcal{U}_n$, $i \in [d_X]$, $a \in A_i$, $j \in [d_Y]$, $b \in B_j$, define the local long-run covariance matrices $\tilde{\Sigma}_{P,n}^{\mathbf{w}^{\varphi(Z)}}(u)$, $\tilde{\Sigma}_{P,n,i,a}^{\mathbf{w}^{\varphi(Z),\varepsilon}}(u)$, $\tilde{\Sigma}_{P,n,j,b}^{\mathbf{w}^{\varphi(Z),\xi}}(u) \in \mathbb{R}^{d_Z \tilde{d}_n \times d_Z \tilde{d}_n}$ for the $\mathbb{R}^{d_Z \tilde{d}_n}$ -valued stationary processes $(\tilde{\mathbf{G}}_n^{\mathbf{w}^{\varphi(Z)}}(u, \mathcal{H}_t^R))_{t \in \mathbb{Z}}$, $(\tilde{\mathbf{G}}_{P,n,i,a}^{\mathbf{w}^{\varphi(Z),\varepsilon}}(u, \mathcal{H}_t^R))_{t \in \mathbb{Z}}$, $(\tilde{\mathbf{G}}_{P,n,j,b}^{\mathbf{w}^{\varphi(Z),\xi}}(u, \mathcal{H}_t^R))_{t \in \mathbb{Z}}$, respectively, by*

$$\begin{aligned}\tilde{\Sigma}_{P,n}^{\mathbf{w}^{\varphi(Z)}}(u) &= \sum_{h \in \mathbb{Z}} \text{Cov}_P(\tilde{\mathbf{G}}_n^{\mathbf{w}^{\varphi(Z)}}(u, \mathcal{H}_0^{\mathbf{w}^{\varphi(Z)}}), \tilde{\mathbf{G}}_n^{\mathbf{w}^{\varphi(Z)}}(u, \mathcal{H}_h^{\mathbf{w}^{\varphi(Z)}})), \\ \tilde{\Sigma}_{P,n,i,a}^{\mathbf{w}^{\varphi(Z),\varepsilon}}(u) &= \sum_{h \in \mathbb{Z}} \text{Cov}_P(\tilde{\mathbf{G}}_{P,n,i,a}^{\mathbf{w}^{\varphi(Z),\varepsilon}}(u, \mathcal{H}_{0,a}^{\mathbf{w}^{\varphi(Z),\varepsilon}}), \tilde{\mathbf{G}}_{P,n,i,a}^{\mathbf{w}^{\varphi(Z),\varepsilon}}(u, \mathcal{H}_{h,a}^{\mathbf{w}^{\varphi(Z),\varepsilon}})), \\ \tilde{\Sigma}_{P,n,j,b}^{\mathbf{w}^{\varphi(Z),\xi}}(u) &= \sum_{h \in \mathbb{Z}} \text{Cov}_P(\tilde{\mathbf{G}}_{P,n,j,b}^{\mathbf{w}^{\varphi(Z),\xi}}(u, \mathcal{H}_{0,b}^{\mathbf{w}^{\varphi(Z),\xi}}), \tilde{\mathbf{G}}_{P,n,j,b}^{\mathbf{w}^{\varphi(Z),\xi}}(u, \mathcal{H}_{h,b}^{\mathbf{w}^{\varphi(Z),\xi}})).\end{aligned}$$

Next, for each $n \in \mathbb{N}$, $P \in \mathcal{P}_n$, $u \in \mathcal{U}_n$, $i \in [d_X]$, $a \in A_i$, $j \in [d_Y]$, $b \in B_j$, define the corresponding integrated long-run covariance matrices $\tilde{\Sigma}_{P,n}^{\mathbf{w}^{\varphi(Z)}}$, $\tilde{\Sigma}_{P,n,i,a}^{\mathbf{w}^{\varphi(Z),\varepsilon}}$, $\tilde{\Sigma}_{P,n,j,b}^{\mathbf{w}^{\varphi(Z),\xi}} \in \mathbb{R}^{d_Z \tilde{c}_n \times d_Z \tilde{c}_n}$ by

$$\begin{aligned}\tilde{\Sigma}_{P,n}^{\mathbf{w}^{\varphi(Z)}} &= \int_{\mathcal{U}_n} \tilde{\Sigma}_{P,n}^{\mathbf{w}^{\varphi(Z)}}(u) \otimes (\phi(u) \phi^\top(u)) du, \\ \tilde{\Sigma}_{P,n,i,a}^{\mathbf{w}^{\varphi(Z),\varepsilon}} &= \int_{\mathcal{U}_n} \tilde{\Sigma}_{P,n,i,a}^{\mathbf{w}^{\varphi(Z),\varepsilon}}(u) \otimes (\phi(u) \phi^\top(u)) du, \\ \tilde{\Sigma}_{P,n,j,b}^{\mathbf{w}^{\varphi(Z),\xi}} &= \int_{\mathcal{U}_n} \tilde{\Sigma}_{P,n,j,b}^{\mathbf{w}^{\varphi(Z),\xi}}(u) \otimes (\phi(u) \phi^\top(u)) du,\end{aligned}$$

where $\phi(u) = (\phi_1(u), \dots, \phi_{\tilde{c}_n}(u))^\top$.

We require the following regularity assumption due to the sieve estimator, which is analogous to Assumption 3.2 from Ding and Zhou [48]. Specifically, for some sequence of collections of distributions $(\mathcal{P}_n)_{n \in \mathbb{N}}$, we impose a distribution-uniform lower bound on the eigenvalues of the integrated long-run covariance matrices.

Assumption 4.7 (Eigenvalue condition for integrated long-run covariance matrices). *Recall $\tilde{\Sigma}_{P,n}^{\mathbf{w}^{\varphi(Z)}}$, $\tilde{\Sigma}_{P,n,i,a}^{\mathbf{w}^{\varphi(Z),\varepsilon}}$, $\tilde{\Sigma}_{P,n,j,b}^{\mathbf{w}^{\varphi(Z),\xi}}$ from Definition 4.3. Assume that there exists a universal constant $\kappa > 0$ such that for all $n \in \mathbb{N}$, $u \in \mathcal{U}_n$, $i \in [d_X]$, $a \in A_i$, $j \in [d_Y]$, $b \in B_j$, we have*

$$\inf_{P \in \mathcal{P}_n} \min(\lambda_{\min}(\tilde{\Sigma}_{P,n}^{\mathbf{w}^{\varphi(Z)}}), \lambda_{\min}(\tilde{\Sigma}_{P,n,i,a}^{\mathbf{w}^{\varphi(Z),\varepsilon}}), \lambda_{\min}(\tilde{\Sigma}_{P,n,j,b}^{\mathbf{w}^{\varphi(Z),\xi}})) \geq \kappa,$$

where $\lambda_{\min}(\cdot)$ is the smallest eigenvalue of the given matrix.

Again, we emphasize that the locally stationary framework in this section fits into the more general triangular array framework from Section 3. Hence, we can use the same cumulative covariance estimator $\hat{Q}_{t,n}^R$ from Subsection 3.5 for the cumulative covariance matrices $Q_{P,t,n}^R = \sum_{s=\mathbb{T}_n^-}^t \Sigma_{P,s,n}^R$, where $\Sigma_{P,s,n}^R = \tilde{\Sigma}_{P,n}^R(s/n)$ denotes the local long-run covariance matrix at time $s \in \mathcal{T}_n$.

4.7 Theoretical result for Sieve-dGCM

The main result of this section is that the Sieve-dGCM test — implemented by running Algorithm 1 with the predictions from the sieve estimator — will have uniformly asymptotic Type-I error control under the previously stated assumptions.

Theorem 4.1. *Suppose that Assumptions 4.1, 4.2, 4.3, 4.4, 4.5, 4.6, 4.7 all hold for the sequence of collections of distributions $(\mathcal{P}_{0,n}^*)_{n \in \mathbb{N}}$, where $\mathcal{P}_{0,n}^* \subset \mathcal{P}_{0,n}^{\text{CI}}$ for each $n \in \mathbb{N}$. Further, suppose that we*

use the sieve time-varying regression estimator from Subsection 4.3 with the basis functions $\{\phi_{\ell_1}(u)\}$, $\{\varphi_{\ell_2}(z)\}$ chosen to be mapped Legendre polynomials, where the numbers of basis functions are chosen to satisfy $\tilde{c}_n = O(\log(T_n))$, $\tilde{d}_n = O(\log(T_n))$. Then the assumptions of Theorem 3.1 hold for $(\mathcal{P}_{0,n}^*)_{n \in \mathbb{N}}$, and the sieve estimators will achieve the convergence rates required by Theorem 3.1.

Throughout this section, we have used Legendre polynomials as the basis functions. In the next section, we investigate the finite sample performance of the Sieve-dGCM test using Legendre basis functions. We emphasize that the Legendre polynomials in our theoretical analysis and simulations can easily be substituted with trigonometric polynomials, wavelets, or other Jacobi polynomials.

5 Numerical Simulations

This section is structured as follows. In Subsection 5.1, we explain how to select the parameters of the sieve estimator and the cumulative covariance estimator. In Subsection 5.2, we report the simulation results. In Subsection 5.3, we discuss the simulation results.

The data generating processes (DGPs) for the simulations below are variants of the DGPs from Shah and Peters [149] with the following changes: (1) the errors are time-varying autoregressive processes instead of iid sequences, (2) the covariates are time-varying autoregressive processes instead of iid sequences, and (3) the regression functions are time-varying instead of time-invariant. We focus our simulations on univariate response settings, as our primary applications of interest are variable selection and causal discovery. We compare our proposed Sieve-dGCM test with the original generalized covariance measure (GCM) test from Shah and Peters [149] and the residual prediction test (RPT) from Shah and Bühlmann [148] and Heinze-Deml, Peters, and Meinshausen [70]. We will report simulation results for other settings and compare with additional conditional independence tests in a separate manuscript.

For the Sieve-dGCM test, we run Algorithm 1 based on predictions from the sieve time-varying nonlinear regression estimator. We use Legendre polynomials as the basis functions as in our theoretical analysis in Section 4. The numbers of basis functions for time and the covariate values were chosen using the subsampling cross-validation procedure from Subsection 5.1. The lag-window parameter for covariance estimation was selected according to the minimum volatility method from Subsection 5.1. We use $s = 5000$ Monte Carlo simulations to approximate the desired quantile. All of the empirical rejection rates below are based on 100 simulated realizations from the nonstationary DGPs.

5.1 Parameter selection via subsampling and minimum volatility

To begin, we introduce a novel cross-validation approach based on subsampling which can be used for selecting the parameters of “global” estimators of time-varying regression functions. The approach we present here is designed for the case where the global estimator is fit once on all the data. When using sequential estimation as in Remark 4.1, standard approaches for time series cross-validation can be used; see Subsection 5.10 of Hyndman [76].

Our approach complements the cross-validation procedure suggested in Subsection 5.1 of Ding and Zhou [48], which is only for parameter selection in the autoregressive forecasting setting. Also, we note that Dahlhaus and Richter [39, 38] theoretically investigated cross-validation for locally stationary processes in the context of selecting bandwidths for kernel smoothing estimators (i.e. a “local” estimation approach). In contrast, our proposed cross-validation approach is for “global” estimators, such as the sieve estimator from Section 4.

The main idea of our cross-validation scheme is to create several folds constructed by sampling the original time series at a *lower sampling frequency*. Specifically, for some buffer $\gamma \in \mathbb{N}_0$ and index $k = 1, \dots, 2(\gamma + 1)$, the k -th fold will consist of the subsampled time series

$$\mathcal{T}_n^{(k)} = \{\mathbb{T}_n^- + k - 1 + 2j(\gamma + 1) : j = 0, 1, \dots, \lfloor \frac{\mathbb{T}_n^+ - \mathbb{T}_n^- - k + 1}{2(\gamma + 1)} \rfloor\}.$$

For instance, when the buffer $\gamma = 0$, we have $\mathcal{T}_n^{(1)} = \{\mathbb{T}_n^-, \mathbb{T}_n^- + 2, \dots\}$ and $\mathcal{T}_n^{(2)} = \{\mathbb{T}_n^- + 1, \mathbb{T}_n^- + 3, \dots\}$. Similarly, when the buffer $\gamma = 1$, we have $\mathcal{T}_n^{(1)} = \{\mathbb{T}_n^-, \mathbb{T}_n^- + 4, \dots\}$, $\mathcal{T}_n^{(2)} = \{\mathbb{T}_n^- + 1, \mathbb{T}_n^- + 5, \dots\}$, $\mathcal{T}_n^{(3)} = \{\mathbb{T}_n^- + 2, \mathbb{T}_n^- + 6, \dots\}$, and $\mathcal{T}_n^{(4)} = \{\mathbb{T}_n^- + 3, \mathbb{T}_n^- + 7, \dots\}$. The reason we refer to γ as a buffer will be made clear below.

We describe our cross-validation scheme in the context of a basic grid search procedure for pedagogical reasons. For each parameter combination, do the following. For each index $k = 1, \dots, \gamma + 1$, use the k -th fold $\mathcal{T}_n^{(k)}$ to estimate the *entire* time-varying regression function (i.e. on a suitably fine grid of rescaled times and covariate values) using the “global” estimator. Afterwards, calculate the residuals based on the observations in the $(k + \gamma + 1)$ -th fold $\mathcal{T}_n^{(k+\gamma+1)}$. By construction, there are γ time points in between the observations in $\mathcal{T}_n^{(k)}$ and $\mathcal{T}_n^{(k+\gamma+1)}$. Next, reverse the roles of the folds. That is, for each index $k = 1, \dots, \gamma + 1$, estimate the *entire* time-varying regression function (i.e. on a suitably fine grid) using the $(k + \gamma + 1)$ -th fold $\mathcal{T}_n^{(k+\gamma+1)}$, and then calculate the corresponding residuals based on the observations in the k -th fold $\mathcal{T}_n^{(k)}$. Finally, for each $k = 1, \dots, 2(\gamma + 1)$, calculate the mean squared error $\text{MSE}^{(k)}$ based on the residuals in fold $\mathcal{T}_n^{(k)}$. Select the parameter combination which yields the lowest average mean squared error

$$\overline{\text{MSE}} = \frac{1}{2(\gamma + 1)} \sum_{k=1}^{2(\gamma+1)} \text{MSE}^{(k)}.$$

In practice, γ should be chosen large enough to account for the temporal dependence, but small enough so that there is enough data to estimate the time-varying regression functions. In our simulations with Sieve-dGCM, we use the buffer $\gamma = 1$ and the grid of parameters $\{1, 2, \dots, 10\} \times \{1, 2, \dots, 10\}$ corresponding to the number of sieve basis functions for time and the covariate values. Note that we allow for the regressions of X on \mathbf{Z} and Y on \mathbf{Z} to have different numbers of basis functions. In future work, we will study the statistical properties of this cross-validation procedure as the buffer $\gamma = \gamma_n$ grows with the sample size n using infill asymptotics. For now, this cross-validation approach serves as a practical technique for parameter selection for generic “global” estimators of time-varying regression functions, such as the sieve estimator.

Next, we discuss how to select the lag-window size parameter L_n for the covariance estimator with a version of the minimum volatility method suggested by Luo and Wu [102]. First, select $H \in \mathbb{N}$ candidate lag-window sizes $l_1 < l_2 < \dots < l_H$. For each index $h = 1, \dots, H$, let

$$\hat{\Sigma}_{t,n,l_h} = \frac{1}{l_h} \left(\sum_{s=t-l_h+1}^t \hat{\mathbf{R}}_{s,n} \right)^{\otimes 2}$$

be the lag-window estimate of the local long-run covariance matrix at time t using the candidate lag-window size $l_h \in \mathbb{N}$. Second, calculate the minimum volatility criterion for each $j = 1, \dots, H$,

$$\mathbf{MV}(j) = \max_{t=\mathbb{T}_n^-+l_H, \dots, \mathbb{T}_n^+} \text{se}[(\hat{\Sigma}_{t,n,l_h})_{h=1 \vee (j-\Delta)}^{H \wedge (j+\Delta)}],$$

where $\Delta \in \mathbb{N}$ is chosen heuristically to balance robustness and adaptivity, and

$$\text{se}[(\hat{\Sigma}_{t,n,l_h})_{h=h_1}^{h_2}] = \text{tr} \left[\frac{1}{h_2 - h_1 + 1} \sum_{h=h_1}^{h_2} \left(\hat{\Sigma}_{t,n,l_h} - \frac{1}{h_2 - h_1 + 1} \sum_{l=h_1}^{h_2} \hat{\Sigma}_{t,n,l_h} \right)^2 \right]^{1/2},$$

with $h_1 = 1 \vee (j - \Delta)$ and $h_2 = H \wedge (j + \Delta)$. Third, select the lag-window size L_n^* that corresponds to the index j^* which yields the smallest minimum volatility criterion

$$j^* = \arg \min_{j=1, \dots, H} \mathbf{MV}(j).$$

We use the following setup in our simulations. We consider $H = \lfloor n/2 \rfloor$ candidate lag-windows with sizes $l_1 = 1, l_2 = 2, \dots, l_H = \lfloor n/2 \rfloor$. We use $\Delta = 12$ so that 25 consecutive lag-window sizes are typically used in the calculation of the minimum volatility criterion $\mathbf{MV}(j)$ for each $j = 1, \dots, H$.

5.2 Analysis of level and power

To begin, we investigate the setting with $d_X = 1$, $d_Y = 1$, $d_Z = 1$ with no time-offsets, so $A = \{0\}$, $B = \{0\}$, $C = \{0\}$, and $\mathcal{T}_n = [n]$. In this case, we can simply refer to the multivariate process as

$$(X_{t,n}, Y_{t,n}, Z_{t,n})_{t \in [n]}.$$

We test for the null hypothesis

$$Y_{t,n} \perp\!\!\!\perp X_{t,n} \mid Z_{t,n} \text{ for all times } t \in [n],$$

versus the alternative hypothesis of

$$Y_{t,n} \not\perp\!\!\!\perp X_{t,n} \mid Z_{t,n} \text{ for all times } t \in [n],$$

because we assume that we can restrict the collection of distributions to be those with time-invariant conditional dependencies. We conduct simulations with the Sieve-dGCM test with $\alpha \in \{0.025, 0.05\}$ for the quantile $\hat{q}_{1-\alpha}^{\text{boot}}$ from Algorithm 1.

For this setting, we use the test statistic based on the maximum absolute value achieved by the partial sum process of residual products. Let

$$Z_{t,n} = 0.4 \left(\frac{2 + \sin(2\pi t/n)}{2} \right) Z_{t-1,n} + \mathcal{N}(0, 1).$$

We use the following DGP for the size simulations with regression complexity parameter $K \in \{1, 2\}$:

$$\begin{aligned} Y_{t,n} &= 0.4 \left(\frac{2 + \sin(2\pi t/n)}{2} \right) \exp(-Z_{t,n}^2) \sin(K Z_{t,n}) + \xi_{t,n}, \\ X_{t,n} &= 0.4 \left(\frac{2 + \sin(2\pi t/n)}{2} \right) \exp(-Z_{t,n}^2) \sin(K Z_{t,n}) + \varepsilon_{t,n}. \end{aligned}$$

We use the following DGP for the power simulations with regression complexity parameter $K \in \{1, 2\}$ and effect size parameter $\beta \in \{0.3, 0.6, 0.9\}$:

$$\begin{aligned} Y_{t,n} &= 0.4 \left(\frac{2 + \sin(2\pi t/n)}{2} \right) \exp(-Z_{t,n}^2) \sin(K Z_{t,n}) + \beta X_{t,n} + \xi_{t,n}, \\ X_{t,n} &= 0.4 \left(\frac{2 + \sin(2\pi t/n)}{2} \right) \exp(-Z_{t,n}^2) \sin(K Z_{t,n}) + \varepsilon_{t,n}. \end{aligned}$$

We use the following error processes for the power and size simulations:

$$\begin{aligned} \varepsilon_{t,n} &= 0.6 \left(\frac{2 + \sin(2\pi t/n)}{2} \right) \varepsilon_{t-1,n} + 0.3 \mathcal{N}(0, 1), \\ \xi_{t,n} &= 0.6 \left(\frac{2 + \sin(2\pi t/n)}{2} \right) \xi_{t-1,n} + 0.3 \mathcal{N}(0, 1). \end{aligned}$$



Figure 2: Empirical rejection rates for our dynamic generalized covariance measure (dGCM) test with the sieve time-varying regression estimator in the $d_X = 1$, $d_Y = 1$, $d_Z = 1$ setting.



Figure 3: Empirical rejection rates for the original generalized covariance measure (GCM) test from [149] with a generalized additive model in the $d_X = 1$, $d_Y = 1$, $d_Z = 1$ setting.

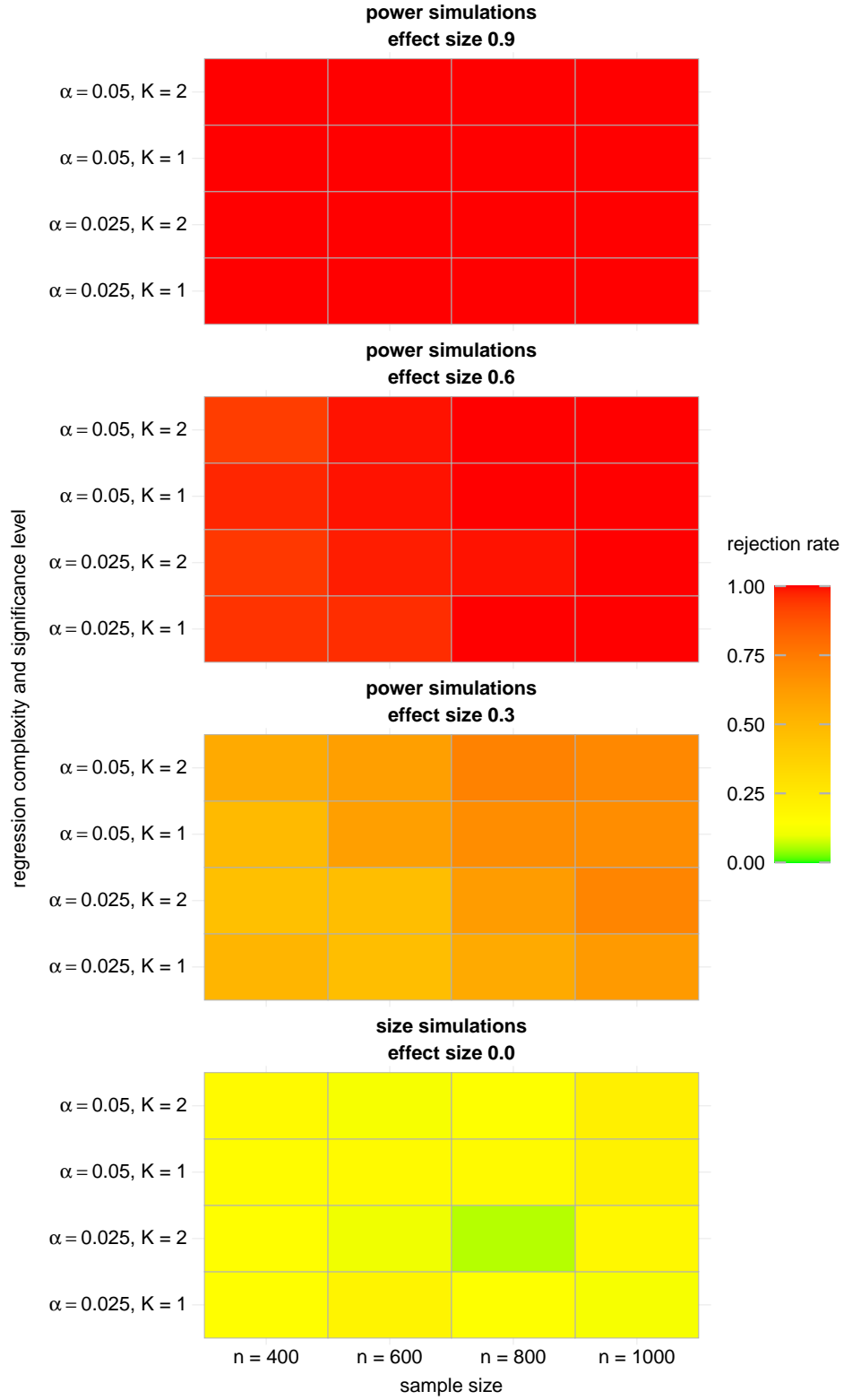


Figure 4: Empirical rejection rates for the residual prediction test (RPT) from [70, 148] with the Nyström method and a random forest model in the $d_X = 1, d_Y = 1, d_Z = 1$ setting.

Next, we consider the setting with $d_X = 1$, $d_Y = 1$, $d_Z = 2$ with no time-offsets, so $A = \{0\}$, $B = \{0\}$, $C = \{0\}$, and $\mathcal{T}_n = [n]$. In this case, we refer to the multivariate process as

$$(X_{t,n}, Y_{t,n}, Z_{t,n,1}, Z_{t,n,2})_{t \in [n]},$$

because there are two covariate processes. Similarly to the previous setting, we test for the null hypothesis of

$$X_{t,n} \perp\!\!\!\perp Y_{t,n} \mid (Z_{t,n,1}, Z_{t,n,2}) \text{ for all times } t \in [n],$$

versus the alternative hypothesis of

$$X_{t,n} \not\perp\!\!\!\perp Y_{t,n} \mid (Z_{t,n,1}, Z_{t,n,2}) \text{ for all times } t \in [n],$$

because we assume that we can restrict the collection of distributions to be those with time-invariant conditional dependencies. We conduct simulations with the Sieve-dGCM test with $\alpha \in \{0.025, 0.05\}$ for the quantile $\hat{q}_{1-\alpha}^{\text{boot}}$ from Algorithm 1.

For this setting, we use the test statistic based on the maximum absolute value achieved by the partial sum process of residual products. Let

$$Z_{t,n,k} = 0.4 \left(\frac{2 + \sin(2\pi t/n)}{2} \right) Z_{t-1,n,k} + \mathcal{N}(0, 1).$$

We use the following DGP for the size simulations with regression complexity parameter $K \in \{1, 2\}$:

$$\begin{aligned} Y_{t,n} &= 0.4 \left(\frac{2 + \sin(2\pi t/n)}{2} \right) \sum_{k=1}^2 \exp(-Z_{t,n,k}^2) \sin(K Z_{t,n,k}) + \xi_{t,n}, \\ X_{t,n} &= 0.4 \left(\frac{2 + \sin(2\pi t/n)}{2} \right) \sum_{k=1}^2 \exp(-Z_{t,n,k}^2) \sin(K Z_{t,n,k}) + \varepsilon_{t,n}. \end{aligned}$$

We use the following DGP for the power simulations with regression complexity parameter $K \in \{1, 2\}$ and effect size parameter $\beta \in \{0.3, 0.6, 0.9\}$:

$$\begin{aligned} Y_{t,n} &= 0.4 \left(\frac{2 + \sin(2\pi t/n)}{2} \right) \sum_{k=1}^2 \exp(-Z_{t,n,k}^2) \sin(K Z_{t,n,k}) + \beta X_{t,n} + \xi_{t,n}, \\ X_{t,n} &= 0.4 \left(\frac{2 + \sin(2\pi t/n)}{2} \right) \sum_{k=1}^2 \exp(-Z_{t,n,k}^2) \sin(K Z_{t,n,k}) + \varepsilon_{t,n}. \end{aligned}$$

We use the following error processes for the power and size simulations:

$$\begin{aligned} \varepsilon_{t,n} &= 0.6 \left(\frac{2 + \sin(2\pi t/n)}{2} \right) \varepsilon_{t-1,n} + 0.3 \mathcal{N}(0, 1), \\ \xi_{t,n} &= 0.6 \left(\frac{2 + \sin(2\pi t/n)}{2} \right) \xi_{t-1,n} + 0.3 \mathcal{N}(0, 1). \end{aligned}$$

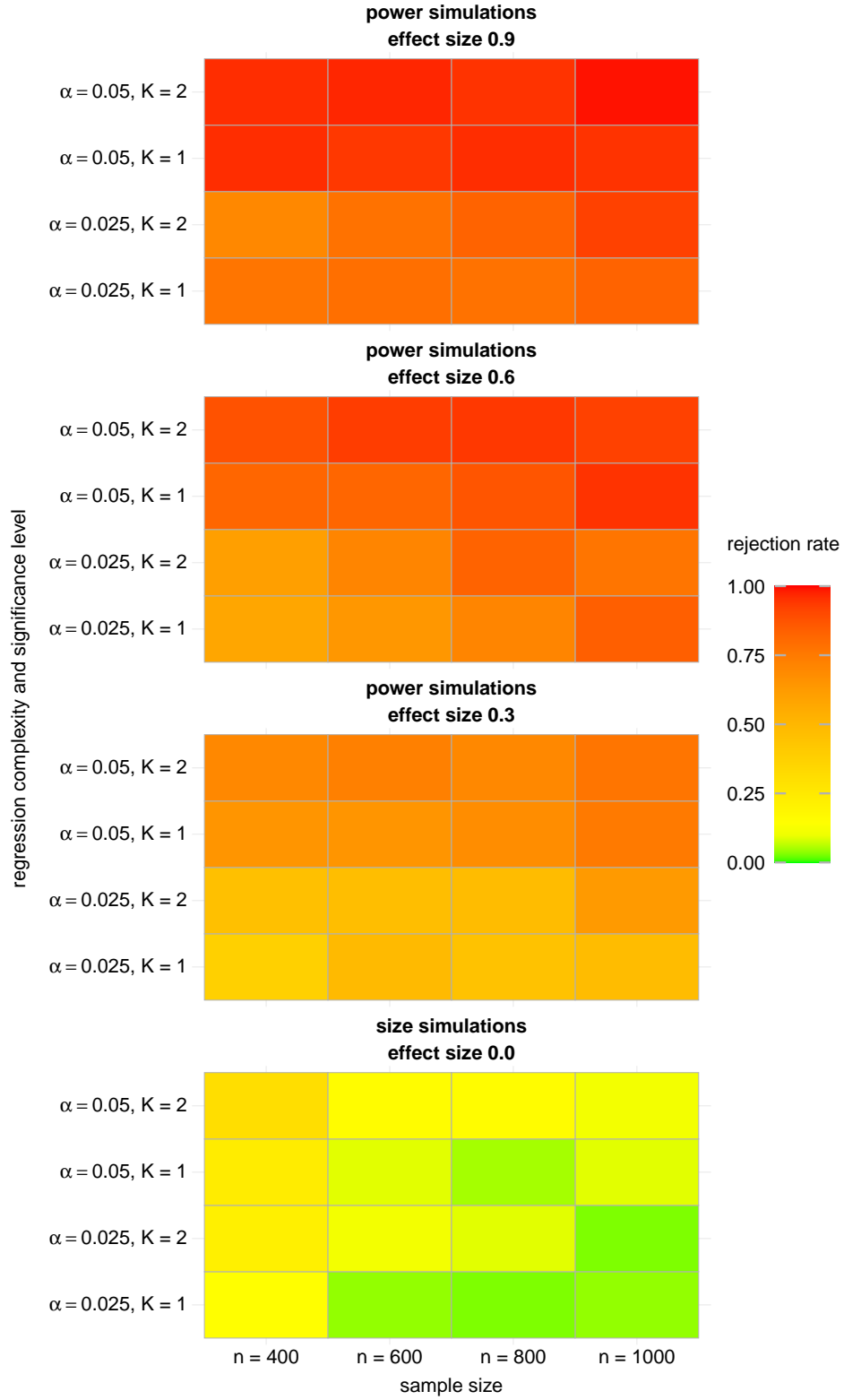


Figure 5: Empirical rejection rates for our dynamic generalized covariance measure (dGCM) test with the sieve time-varying regression estimator in the $d_X = 1$, $d_Y = 1$, $d_Z = 2$ setting.



Figure 6: Empirical rejection rates for the original generalized covariance measure (GCM) test from [149] with a generalized additive model in the $d_X = 1$, $d_Y = 1$, $d_Z = 2$ setting.

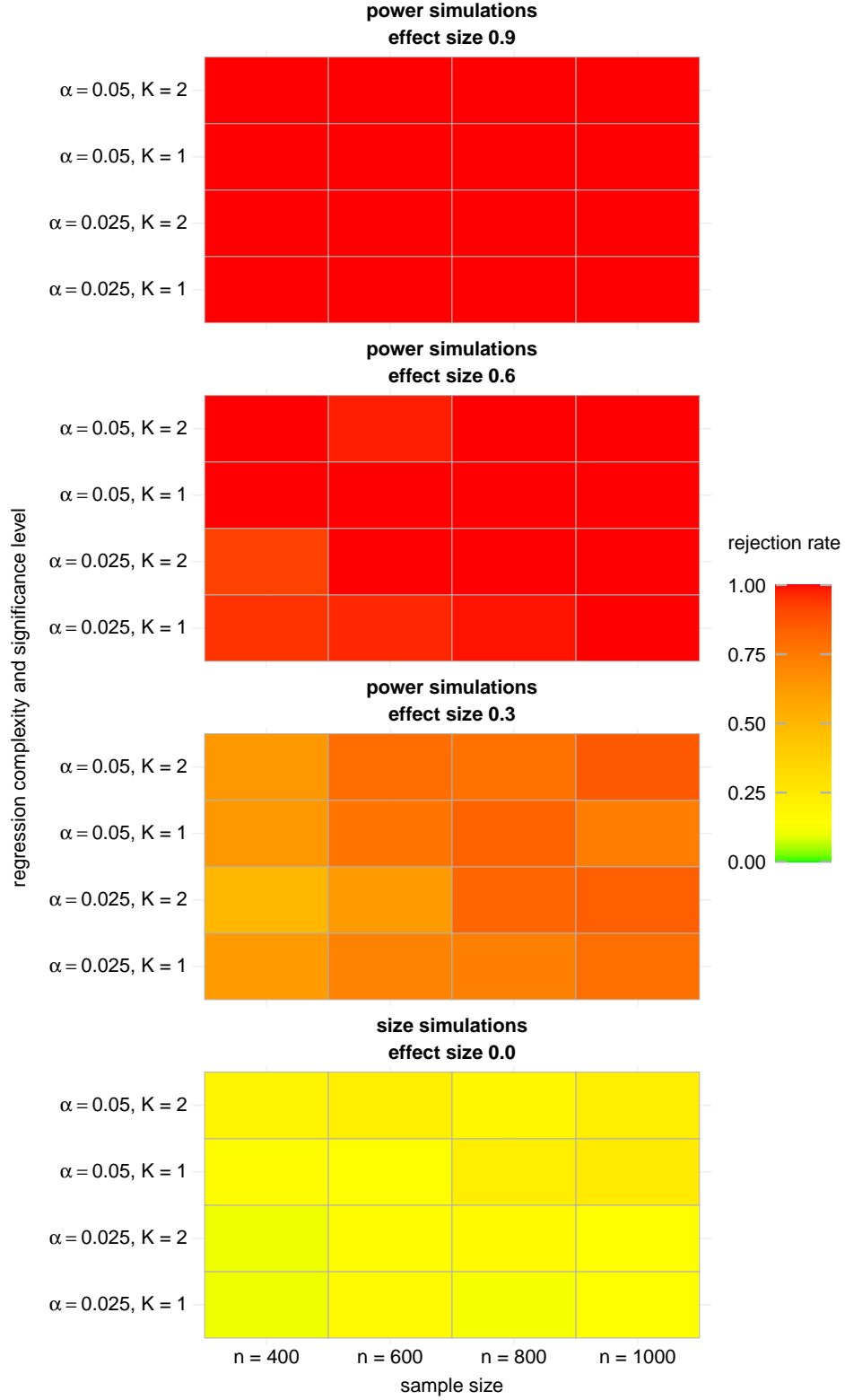


Figure 7: Empirical rejection rates for the residual prediction test (RPT) from [70, 148] with the Nyström method and a random forest model in the $d_X = 1, d_Y = 1, d_Z = 2$ setting.

5.3 Discussion

The main takeaways are as follows. All of the conditional independence tests are able to detect conditional dependencies. However, GCM and RPT fail to hold the level. In contrast, we find that dGCM can hold the level if the sample size is large enough to reliably estimate the time-varying regression functions. This is not entirely surprising, as RPT and GCM were developed for the idealized iid setting and do not account for autocorrelation or nonstationarity. Nevertheless, it is useful to see how these tests perform when relatively mild forms of nonstationarity and autocorrelation are added to the DGPs from Shah and Peters [149]. We conclude that the dGCM test shows promise for reliably detecting conditional dependencies in non-iid settings. Hence, we regard the dGCM test as a nonstationarity-and-autocorrelation robust extension of the GCM test from Shah and Peters [149], thereby contributing to the growing collection of GCM-inspired tests [144, 100, 33, 164, 78, 25, 101].

In future work, we plan to exhaustively compare the dGCM test with a large number of existing conditional independence tests across a diverse range of non-iid settings and challenging tasks. Examples include detecting time-delayed causal effects, screening out irrelevant forecasting signals, handling high-dimensional covariates and responses, and exploring more extreme forms of autocorrelation and nonstationarity. Also, we plan to investigate the performance of the dGCM test with other time-varying regression estimators.

We emphasize that our simulations adhere to the asymptotic framework of locally stationary processes as described in Section 4. Within this framework, as n grows we gain more and more observations of each *local structure* of the nonstationary process. Standard long-run asymptotics for stationary processes emerge as a special case of this infill asymptotic framework when all parameter curves and regression functions are time-invariant. Likewise, the iid setting arises as another special case when there is neither nonstationarity nor temporal dependence.

We provide an empirical demonstration of the “no-free-lunch in conditional independence testing” results from [149, 19] with our simulations by increasing the complexity of the regression functions while holding the sample size fixed. When the regression complexity (parametrized by K) is large relative to the sample size n , we observe a degradation in Type-I error control because the regression functions cannot be reliably estimated. This phenomena can also be observed in the simulation experiments with the GCM test; see the discussion in Section 5 of Shah and Peters [149]. In light of the no-free-lunch results, we understand that this is inevitable: no matter how large the sample size, it is impossible to ensure the correct significance level for every null distribution. Hence, there will always be some null distribution in which the Type-I error rate exceeds the prespecified significance level.

Crucially, the uniform level guarantee for our test only applies when the assumptions on the *sequence* of collections of distributions are satisfied. In the context of our numerical simulations, this means that any sequence of distributions parametrized by a sequence $(K_n)_{n \in \mathbb{N}}$ in which the regression complexity parameter K_n grows with the sample size n at some rate will violate Assumption 4.6. Therefore, the uniform asymptotic Type-I error control guarantee we provide is not applicable. On the other hand, if we increase the sample size while holding the degree of complexity of the regression functions fixed (e.g. setting $K_n = K = 1$ or $K_n = K = 2$ for all sample sizes n), then the uniform asymptotic Type-I error control guarantee is applicable. This explains why conditional independence tests may fail to control Type-I error and highlights the “honesty” of uniform level guarantees [94, 82, 160, 129, 88].

5.4 Accuracy of the strong Gaussian approximation

Lastly, we consider the hypothetical scenario in which the time-varying regression functions are perfectly estimated. We conduct these oracle simulations with the dGCM test with $\alpha \in \{0.025, 0.05\}$ for the quantile $\hat{q}_{1-\alpha}^{\text{boot}}$ from Algorithm 1. For the power simulations, the residual products will be $\varepsilon_{t,n}(\beta\varepsilon_{t,n} + \xi_{t,n})$ with effect sizes $\beta \in \{0.3, 0.6, 0.9\}$. For the size simulations, the residual products will just be the error products $\varepsilon_{t,n}\xi_{t,n}$. We use the same error processes from the previous simulations:

$$\begin{aligned}\varepsilon_{t,n} &= 0.6 \left(\frac{2 + \sin(2\pi t/n)}{2} \right) \varepsilon_{t-1,n} + 0.3\mathcal{N}(0, 1), \\ \xi_{t,n} &= 0.6 \left(\frac{2 + \sin(2\pi t/n)}{2} \right) \xi_{t-1,n} + 0.3\mathcal{N}(0, 1).\end{aligned}$$

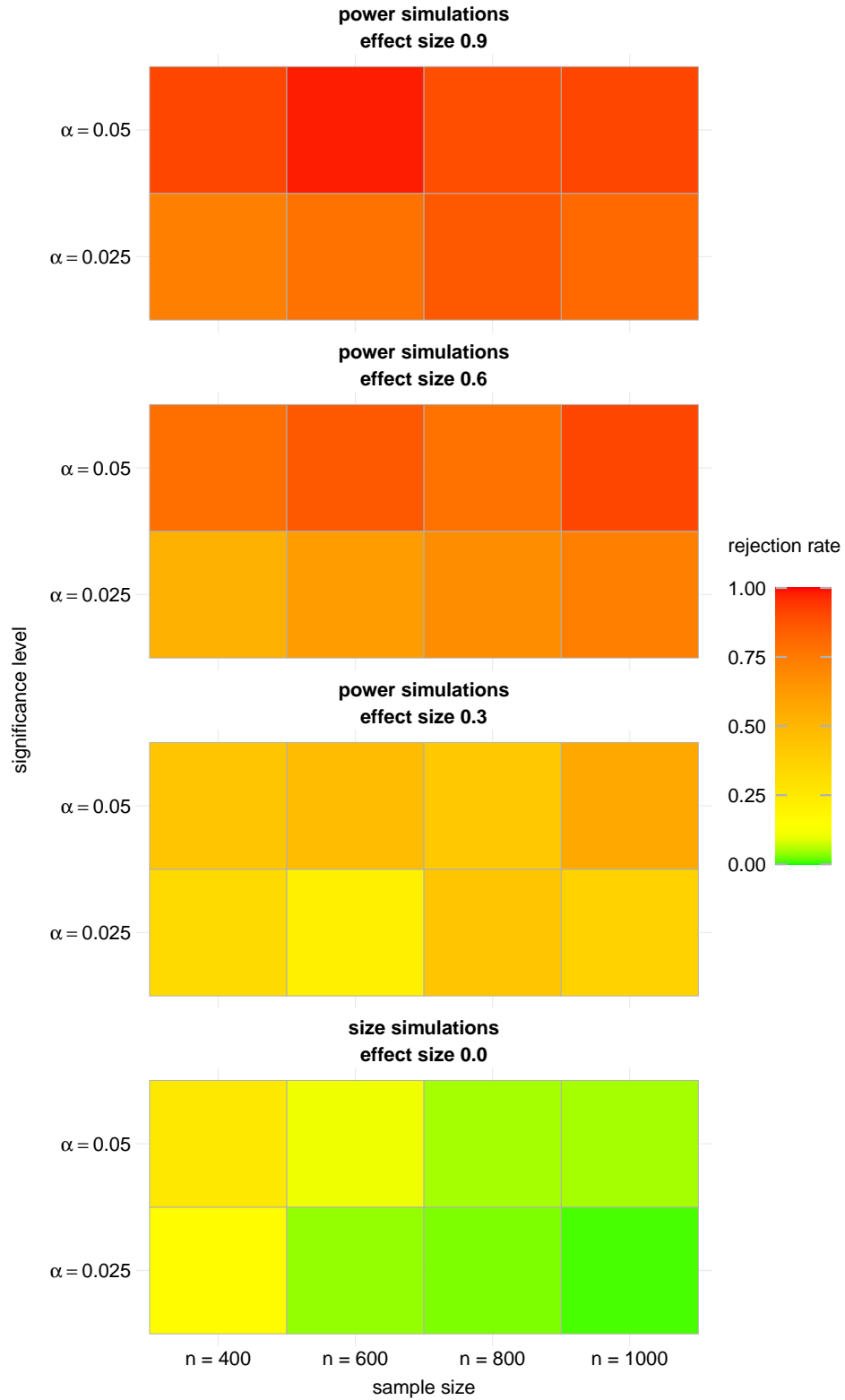


Figure 8: Empirical rejection rates for our dynamic generalized covariance measure (dGCM) test with perfectly estimated time-varying regression functions.

6 Future Work

In this paper, we introduced a nonstationarity-and-autocorrelation robust test for conditional independence. In our companion paper [166], we apply the dGCM test to nonstationary time series data from the domain of epidemiology. Next, we discuss three promising avenues for exciting future work.

First, we plan to develop statistical techniques for nonstationary nonlinear time series which utilize our conditional independence test as a core component. For example, a Markov blanket-based variable selection procedure for forecasting and a causal discovery algorithm for nonstationary nonlinear time series. We note that it may be possible to use the theoretical tools for nonlinear locally stationary processes and the functional dependence measure to develop a unified causal inference framework for nonstationary processes by building on prior work for stationary processes [143, 127, 139, 140, 137].

Second, we plan to explore topics related to time-varying regression estimation. It would be of interest to theoretically investigate the sieve estimator from Section 4 in the high-dimensional setting by introducing a sparsity-inducing penalty for regularization as in Zhang and Simon [185]. It would also be of interest to develop a computationally efficient online estimation procedure for the sieve estimator by taking inspiration from Zhang and Simon [184]. Along the way, we plan to theoretically investigate our subsampling cross-validation procedure from Subsection 5.1 for selecting the parameters of global estimators of time-varying regression functions. Additionally, it would be of interest to develop guarantees for time-varying nonlinear regression estimators in the context of processes with the total variation-type nonstationarity condition from Assumption 3.6. It may also be possible to investigate the convergence rates for deep neural network regression estimators as in Kurisu, Fukami, and Koike [89], but in the context of nonstationary processes with the functional dependence measure.

Third, there are several possible future research directions for conditional independence testing in this setting. While our test is based on the expected conditional covariance functional, our framework can easily be adapted to use any other functional equal to zero under the null of conditional independence. In particular, using higher-order functionals may be of interest in more complicated settings because the expected conditional covariance functional lacks sensitivity to nonlinear relationships and interactions; see Zhang and Janson [183] for more discussion. Specifically, it would be valuable to develop such tests without compromising on practicality, which is one of the key advantages of our regression-based approach. In Subsection A.5, we discuss various conditional independence tests designed specifically for the locally stationary setting. However, those test statistics utilize kernel smoothing, so the resulting tests can be very sensitive to the choice of the bandwidth parameters.

References

- [1] Alekh Agarwal and John C. Duchi. “The generalization ability of online algorithms for dependent data”. *IEEE Transactions on Information Theory* 59.1 (2012), pp. 573–587.
- [2] Pierre Alquier, Xiaoyin Li, and Olivier Wintenberger. “Prediction of time series by statistical learning: general losses and fast rates”. *Dependence Modeling* 1 (2013), pp. 65–93.
- [3] Pierre-Olivier Amblard and Olivier J. J. Michel. “The relation between Granger causality and directed information theory: a review”. *Entropy* 15.1 (2012), pp. 113–143.
- [4] Donald W.K. Andrews and Panle Jia Barwick. “Inference for parameters defined by moment inequalities: a recommended moment selection procedure”. *Econometrica: Journal of the Econometric Society* 80.6 (2012), pp. 2805–2826.
- [5] Donald W.K. Andrews and Patrik Guggenberger. “Validity of subsampling and plug-in asymptotic inference for parameters defined by moment inequalities”. *Econometric Theory* 25.3 (2009), pp. 669–709.
- [6] Donald W.K. Andrews and Gustavo Soares. “Inference for parameters defined by moment inequalities using generalized moment selection”. *Econometrica: Journal of the Econometric Society* 78.1 (2010), pp. 119–157.
- [7] Lujia Bai and Weichi Wu. “Time-varying correlation network analysis of non-stationary multivariate time series with complex trends”. arXiv preprint arXiv:2302.05158. 2023.
- [8] Rina Foygel Barber, Emmanuel J. Candès, and Richard J. Samworth. “Robust inference with knockoffs”. *The Annals of Statistics* 48.3 (2020), pp. 1409–1431.

- [9] Rina Foygel Barber and Lucas Janson. “Testing goodness-of-fit and conditional independence with approximate co-sufficient sampling”. *The Annals of Statistics* 50.5 (2022), pp. 2514–2544.
- [10] Sumanta Basu and George Michailidis. “Regularized estimation in sparse high-dimensional time series models”. *Annals of Statistics* 43.4 (2015), pp. 1535–1567.
- [11] Sumanta Basu and Suhasini Subba Rao. “Graphical models for nonstationary time series”. *The Annals of Statistics* 51.4 (2023), pp. 1453–1483.
- [12] Sumanta Basu, Ali Shojaie, and George Michailidis. “Network Granger causality with inherent grouping structure”. *The Journal of Machine Learning Research* 16.1 (2015), pp. 417–453.
- [13] Carina Beering. “A functional central limit theorem and its bootstrap analogue for locally stationary processes with application to independence testing”. PhD Dissertation, Technische Universität Braunschweig, 2021.
- [14] Alexandre Belloni et al. “Some new asymptotic theory for least squares series: pointwise and uniform results”. *Journal of Econometrics* 186.2 (2015), pp. 345–366.
- [15] Yoav Benjamini and Yosef Hochberg. “Controlling the false discovery rate: a practical and powerful approach to multiple testing”. *Journal of the Royal Statistical Society Series B: Statistical Methodology* 57.1 (1995), pp. 289–300.
- [16] Yoav Benjamini and Daniel Yekutieli. “The control of the false discovery rate in multiple testing under dependency”. *Annals of Statistics* 29.4 (2001), pp. 1165–1188.
- [17] W.R. Bennett. “Statistics of regenerative digital transmission”. *Bell System Technical Journal* 37.6 (1958), pp. 1501–1542.
- [18] Peter Bloomfield, Harry L. Hurd, and Robert B. Lund. “Periodic correlation in stratospheric ozone data”. *Journal of Time Series Analysis* 15.2 (1994), pp. 127–150.
- [19] Juraĳ Bodik and Olivier C. Pasche. “Granger causality in extremes”. arXiv preprint arXiv:2407.09632. 2024.
- [20] Soham Bonnerjee, Sayar Karmakar, and Wei Biao Wu. “Gaussian approximation for non-stationary time series with optimal rate and explicit construction”. *The Annals of Statistics* 52.5 (2024), pp. 2293–2317.
- [21] Taoufik Bouezmarni, Jeroen V.K. Rombouts, and Abderrahim Taamouti. “Nonparametric copula-based test for conditional independence with applications to Granger causality”. *Journal of Business and Economic Statistics* 30.2 (2012), pp. 275–287.
- [22] Leo Breiman. “Random forests”. *Machine Learning* 45 (2001), pp. 5–32.
- [23] Guy-Niklas Brunotte. “A test of independence under local stationarity based on the local characteristic function”. 2022.
- [24] Emmanuel J. Candés et al. “Panning for gold: model-X knockoffs for high dimensional controlled variable selection”. *Journal of the Royal Statistical Society Series B: Statistical Methodology* 80.3 (2018), pp. 551–577.
- [25] Abhinav Chakraborty, Jeffrey Zhang, and Eugene Katsevich. “Doubly robust and computationally efficient high-dimensional variable selection” (2024). arXiv preprint arXiv:2409.09512.
- [26] Jinyuan Chang, Xiaohui Chen, and Mingcong Wu. “Central limit theorems for high dimensional dependent data”. *Bernoulli* 30.1 (2024), pp. 712–742.
- [27] Likai Chen, Ekaterina Smetanina, and Wei Biao Wu. “Estimation of nonstationary nonparametric regression model with multiplicative structure”. *The Econometrics Journal* 25.1 (2022), pp. 176–214.
- [28] Tianqi Chen and Carlos Guestrin. “Xgboost: a scalable tree boosting system”. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2016), pp. 785–794.
- [29] Xiaohong Chen. “Large sample sieve estimation of semi-nonparametric models”. *Handbook of Econometrics* 6 (2007), pp. 5549–5632.
- [30] Victor Chernozhukov, Denis Chetverikov, and Kengo Kato. “Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors”. *The Annals of Statistics* 41.6 (2013), pp. 2786–2819.

- [31] Victor Chernozhukov, Denis Chetverikov, and Kengo Kato. “Inference on causal and structural parameters using many moment inequalities”. *The Review of Economic Studies* 86.5 (2019), pp. 1867–1900.
- [32] Victor Chernozhukov, Denis Chetverikov, and Kengo Kato. “Testing many moment inequalities”. 2016.
- [33] Alexander Mangulad Christgau, Lasse Petersen, and Niels Richard Hansen. “Nonparametric conditional local independence testing”. *The Annals of Statistics* 51.5 (2022), pp. 2116–2144.
- [34] Kyle Cranmer, Johann Brehmer, and Gilles Louppe. “The frontier of simulation-based inference”. *Proceedings of the National Academy of Sciences* 117.48 (2020), pp. 30055–30062.
- [35] Yan Cui and Zhou Zhou. “Optimal Short-Term Forecast for Locally Stationary Functional Time Series”. *IEEE Transactions on Information Theory* (2025).
- [36] Rainer Dahlhaus. “Fitting time series models to nonstationary processes”. *The Annals of Statistics* 25.1 (1997), pp. 1–37.
- [37] Rainer Dahlhaus. “Locally stationary processes”. *Handbook of statistics* 30 (2012), pp. 351–413.
- [38] Rainer Dahlhaus and Stefan Richter. “Adaptation for nonparametric estimators of locally stationary processes”. *Econometric Theory* 39.6 (2023), pp. 1123–1153.
- [39] Rainer Dahlhaus and Stefan Richter. “Cross validation for locally stationary processes”. *Annals of Statistics* 47.4 (2019), pp. 2145–2173.
- [40] Rainer Dahlhaus, Stefan Richter, and Wei Biao Wu. “Towards a general theory for nonlinear locally stationary processes”. *Bernoulli* 25.2 (2019), pp. 1013–1044.
- [41] J. J. Daudin. “Partial association measures and an application to qualitative regression”. *Biometrika* 67.3 (1980), pp. 581–590.
- [42] Richard A. Davis and Mikkel S. Nielsen. “Modeling of time series using random forests: theoretical developments”. *Electronic Journal of Statistics* (2020), pp. 3644–3671.
- [43] Philip A. Dawid and Ambuj Tewari. “On learnability under general stochastic processes”. *Harvard Data Science Review* 4.4 (2022).
- [44] Holger Dette, Weichi Wu, and Zhou Zhou. “Change point analysis of correlation in non-stationary time series”. *Statistica Sinica* 29.2 (2019), pp. 611–643.
- [45] Xiucan Ding and Zhou Zhou. “Autoregressive approximations to nonstationary time series with inference and applications”. *The Annals of Statistics* 51.3 (2023), pp. 1207–1231.
- [46] Xiucan Ding and Zhou Zhou. “Estimation and inference for precision matrices of nonstationary time series”. *Annals of Statistics* 48.4 (2020), pp. 2455–2477.
- [47] Xiucan Ding and Zhou Zhou. “On the partial autocorrelation function for locally stationary time series: characterization, estimation and inference”. arXiv preprint arXiv:2401.15778. 2024.
- [48] Xiucan Ding and Zhou Zhou. “Simultaneous sieve inference for time-inhomogeneous nonlinear time series regression”. arXiv preprint arXiv:2112.08545. 2021.
- [49] Xinshuai Dong et al. “On the three demons in causality in finance: time resolution, nonstationarity, and latent factors”. arXiv preprint arXiv:2401.05414. 2023.
- [50] Gary Doran et al. “A permutation-based kernel conditional independence test”. *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence* (2014), pp. 132–141.
- [51] Michael Eichler. “Graphical modelling of multivariate time series”. *Probability Theory and Related Fields* 153 (2012), pp. 233–268.
- [52] Jianqing Fan, Yang Feng, and Lucy Xia. “A projection-based conditional dependence measure with applications to high-dimensional undirected graphical models”. *Journal of Econometrics* 218.1 (2020), pp. 119–139.
- [53] Muhammad Hasan Ferdous, Uzma Hasan, and Md Osman Gani. “Cdans: temporal causal discovery from autocorrelated and non-stationary time series data”. *Proceedings of Machine Learning Research* 219 (2023), pp. 186–207.

- [54] Seth R. Flaxman, Daniel B. Neill, and Alexander J. Smola. “Gaussian processes for independence tests with non-iid data in causal inference”. *ACM Transactions on Intelligent Systems and Technology* 7.2 (2015), pp. 1–23.
- [55] Jean-Pierre Florens and Michel Mouchart. “A note on noncausality”. *Econometrica: Journal of the Econometric Society* (1982), pp. 583–591.
- [56] David T. Frazier and Bonsoo Koo. “Indirect inference for locally stationary models”. *Journal of Econometrics* 223.1 (2021), pp. 1–27.
- [57] Kenji Fukumizu et al. “Kernel measures of conditional dependence”. *Advances in Neural Information Processing Systems* 20 (2007), pp. 489–496.
- [58] William A. Gardner. *Cyclostationarity in communications and signal processing*. IEEE Press, 1994.
- [59] William A. Gardner, Antonio Napolitano, and Luigi Paura. “Cyclostationarity: half a century of research”. *Signal Processing* 86.4 (2006), pp. 639–697.
- [60] Benjamin Goehry. “Random forests for time-dependent processes”. *ESAIM: Probability and Statistics* 24 (2020), pp. 801–826.
- [61] Christian Gourieroux, Alain Monfort, and Eric Renault. “Indirect inference”. *Journal of Applied Econometrics* 8 (1993), S85–S118.
- [62] Clive WJ Granger. “Investigating causal relations by econometric models and cross-spectral methods”. *Econometrica: Journal of the Econometric Society* 37.3 (1969), pp. 424–438.
- [63] Clive WJ Granger. “Testing for causality: a personal viewpoint”. *Journal of Economic Dynamics and Control* 2 (1980), pp. 329–352.
- [64] Arthur Gretton et al. “A kernel statistical test of independence”. *Advances in Neural Information Processing Systems* 20 (2007).
- [65] Peter Grünwald, Alexander Henzi, and Tyron Lardy. “Anytime-valid tests of conditional independence under model-X”. *Journal of the American Statistical Association* 119.546 (2024), pp. 1554–1565.
- [66] Shuva Gupta. “A note on the asymptotic distribution of LASSO estimator for correlated data”. *Sankhya A* 74.1 (2012), pp. 10–28.
- [67] Hanyuan Hang and Ingo Steinwart. “Fast learning from alpha-mixing observations”. *Journal of Multivariate Analysis* 127 (2014), pp. 184–199.
- [68] Steve Hanneke. “Learning whenever learning is possible: universal learning under general stochastic processes”. *The Journal of Machine Learning Research* 22.130 (2021), pp. 1–116.
- [69] Steve Hanneke and Liu Yang. “Statistical learning under nonstationary mixing processes”. *International Conference on Artificial Intelligence and Statistics* 22.1 (2019), pp. 1678–1686.
- [70] Christina Heinze-Deml, Jonas Peters, and Nicolai Meinshausen. “Invariant causal prediction for nonlinear models”. *Journal of Causal Inference* 6.2 (2018), pp. 6887–6909.
- [71] Oliver Hines et al. “Demystifying statistical learning based on efficient influence functions”. *The American Statistician* 76.3 (2022), pp. 292–304.
- [72] Biwei Huang et al. “Causal discovery from heterogeneous/nonstationary data”. *Journal of Machine Learning Research* 21.89 (2020), pp. 1–53.
- [73] Dongming Huang and Lucas Janson. “Relaxing the assumptions of knockoffs by conditioning”. *The Annals of Statistics* 48.5 (2020), pp. 3021–3042.
- [74] Jianhua Z. Huang. “Projection estimation in multiple regression with application to functional ANOVA models”. *The Annals of Statistics* 26.1 (1998), pp. 242–272.
- [75] Tzee-Ming Huang. “Testing conditional independence using maximal nonlinear conditional correlation”. *The Annals of Statistics* 38.4 (2010), pp. 2047–2091.
- [76] R. J. Hyndman. *Forecasting: principles and practice*. OTexts, 2018.
- [77] Guido W. Imbens and Charles F. Manski. “Confidence intervals for partially identified parameters”. *Econometrica: Journal of the Econometric Society* 72.6 (2004), pp. 1845–1857.

- [78] Iden Kalemaj, Shiva Kasiviswanathan, and Aaditya Ramdas. “Differentially private conditional independence testing”. *International Conference on Artificial Intelligence and Statistics* 238 (2024), pp. 3700–3708.
- [79] Olav Kallenberg. *Foundations of modern probability*. Third edition. Springer, 2021.
- [80] Rajeeva L. Karandikar and Mathukumalli Vidyasagar. “Rates of uniform convergence of empirical means with mixing processes”. *Statistics and Probability Letters* 58.3 (2002), pp. 297–307.
- [81] Sayar Karmakar and Wei Biao Wu. “Optimal Gaussian approximation for multiple time series”. 30.3 (2020), pp. 1399–1417.
- [82] Maximilian Kasy. “Uniformity and the delta method”. *Journal of Econometric Methods* 8.1 (2018).
- [83] Guolin Ke et al. “Lightgbm: a highly efficient gradient boosting decision tree”. *Advances in Neural Information Processing Systems* 30 (2017).
- [84] Edward H. Kennedy. “Semiparametric doubly robust targeted double machine learning: a review”. *Handbook of Statistical Methods for Precision Medicine*. Chapman and Hall/CRC, 2024, pp. 207–236.
- [85] Ilmun Kim et al. “Local permutation tests for conditional independence”. *The Annals of Statistics* 50.6 (2022), pp. 3388–3414.
- [86] Jonas Krampe and Suhasini Subba Rao. “Inverse covariance operators of multivariate nonstationary time series”. *Bernoulli* 30.2 (2024), pp. 1177–1196.
- [87] Arun K. Kuchibhotla, John E. Kolassa, and Todd A. Kuffner. “Post-selection inference”. *Annual Review of Statistics and Its Application* 9.1 (2022), pp. 505–527.
- [88] Arun Kumar Kuchibhotla, Sivaraman Balakrishnan, and Larry Wasserman. “Median regularity and honest inference”. *Biometrika* 110.3 (2023), pp. 831–838.
- [89] Daisuke Kurisu, Riku Fukami, and Yuta Koike. “Adaptive deep learning for nonlinear time series models”. *Bernoulli* 31.1 (2025), pp. 240–270.
- [90] Vitaly Kuznetsov and Mehryar Mohri. “Generalization bounds for non-stationary mixing processes”. *Machine Learning* 106.1 (2017), pp. 93–117.
- [91] Vitaly Kuznetsov and Mehryar Mohri. “Generalization bounds for time series prediction with non-stationary processes”. *Algorithmic Learning Theory* 25 (2014), pp. 260–274.
- [92] Vitaly Kuznetsov and Mehryar Mohri. “Learning theory and algorithms for forecasting non-stationary time series”. *Advances in Neural Information Processing Systems* 28 (2015).
- [93] Jia Li, Zhipeng Liao, and Wenyu Zhou. “A general test for functional inequalities”. 2022.
- [94] Ker-Chau Li. “Honest confidence regions for nonparametric regression”. *The Annals of Statistics* 17.3 (1989), pp. 1001–1008.
- [95] Lingling Li et al. “Higher order inference on a treatment effect under low regularity conditions”. *Statistics and Probability Letters* 81.7 (2011), pp. 821–828.
- [96] Molei Liu et al. “Fast and powerful conditional randomization testing via distillation”. *Biometrika* 109.2 (2022), pp. 277–293.
- [97] Weidong Liu, Han Xiao, and Wei Biao Wu. “Probability and moment inequalities under dependence”. *Statistica Sinica* 23.3 (2013), pp. 1257–1272.
- [98] Zhaolu Liu et al. “Kernel-based joint independence tests for multivariate stationary and non-stationary time series”. *Royal Society Open Science* 10.11 (2023).
- [99] Ignacio N. Lobato. “Testing that a dependent process is uncorrelated”. *Journal of the American Statistical Association* 96.455 (2001), pp. 1066–1076.
- [100] Anton Rask Lundborg, Rajen D. Shah, and Jonas Peters. “Conditional independence testing in Hilbert spaces with applications to functional data analysis”. *Journal of the Royal Statistical Society Series B: Statistical Methodology* 84.5 (2022), pp. 1821–1850.
- [101] Anton Rask Lundborg et al. “The projected covariance measure for assumption-lean variable significance testing”. *The Annals of Statistics* 52.6 (2024), pp. 2851–2878.

- [102] Tianpai Luo and Weichi Wu. “Simultaneous inference for monotone and smoothly time varying functions under complex temporal dynamics”. arXiv preprint arXiv:2310.02177. 2023.
- [103] Helmut Lütkepohl. *New introduction to multiple time series analysis*. Springer Science and Business Media, 2005.
- [104] Daniel Malinsky and Peter Spirtes. “Learning the structure of a nonstationary vector autoregression”. *International Conference on Artificial Intelligence and Statistics* 89 (2019), pp. 2986–2994.
- [105] Georg Manten et al. “Signature kernel conditional independence tests in causal discovery for stochastic processes”. arXiv preprint arXiv:2402.18477. 2024.
- [106] Dimitris Margaritis. “Distribution-free learning of Bayesian network structure in continuous domains”. *AAAI* 5 (2005), pp. 825–830.
- [107] Daniel McFadden. “A method of simulated moments for estimation of discrete response models without numerical integration”. *Econometrica: Journal of the Econometric Society* 57.5 (1989), pp. 995–1026.
- [108] Fabian Mies. “Strong Gaussian approximations with random multipliers”. arXiv preprint arXiv:2412.14346. 2024.
- [109] Fabian Mies and Ansgar Steland. “Sequential Gaussian approximation for nonstationary time series in high dimensions”. *Bernoulli* 29.4 (2023), pp. 3114–3140.
- [110] Mehryar Mohri and Vitaly Kuznetsov. “Discrepancy-based theory and algorithms for forecasting non-stationary time series”. *Annals of Mathematics and Artificial Intelligence* 88.4 (2020), pp. 367–399.
- [111] Antonio Napolitano. “Cyclostationarity: limits and generalizations”. *Signal Processing* 120 (2016), pp. 323–347.
- [112] Antonio Napolitano. “Cyclostationarity: new trends and applications”. *Signal Processing* 120 (2016), pp. 385–408.
- [113] Whitney K. Newey. “Convergence rates and asymptotic normality for series estimators”. *Journal of Econometrics* 79.1 (1997), pp. 147–168.
- [114] Whitney K. Newey and James R. Robins. “Cross-fitting and fast remainder rates for semiparametric estimation”. arXiv preprint arXiv:1801.09138. 2018.
- [115] Matey Neykov, Sivaraman Balakrishnan, and Larry Wasserman. “Minimax optimal conditional independence testing”. *The Annals of Statistics* 49.4 (2021), pp. 2151–2177.
- [116] Ziang Niu et al. “Reconciling model-X and doubly robust approaches to conditional independence testing”. *The Annals of Statistics* 52.3 (2024), pp. 895–921.
- [117] Emanuel Parzen and Marcello Pagano. “An approach to modeling seasonally stationary time series”. *Journal of Econometrics* 9.1-2 (1979), pp. 137–153.
- [118] Hoyer Patrik et al. “Nonlinear causal discovery with additive noise models”. *Advances in Neural Information Processing Systems* 21 (2009).
- [119] Judea Pearl. *Causality*. Cambridge University Press, 2009.
- [120] Judea Pearl. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Elsevier, 2014.
- [121] Ling Peng, Yan Zhu, and Wenxuan Zhong. “Lasso regression in sparse linear model with phmixing errors”. *Metrika* 86.1 (2023), pp. 1–26.
- [122] Jonas Peters et al. “Causal discovery with continuous additive noise models”. *Journal of Machine Learning Research* 15.58 (2014), pp. 2009–2053.
- [123] Gilles Pisier. *Martingales in Banach spaces*. Vol. 155. Cambridge University Press, 2016.
- [124] Christopher J. Quinn, Negar Kiyavash, and Todd P. Coleman. “Directed information graphs”. *IEEE Transactions on Information Theory* 61.12 (2015), pp. 6887–6909.
- [125] Joseph D. Ramsey. “A scalable conditional independence test for nonlinear, non-gaussian data”. arXiv preprint arXiv:1401.5031. 2014.

- [126] Rolando Rebolledo. “Central limit theorems for local martingales”. *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete* 51.3 (1980), pp. 269–286.
- [127] Nicolas-Domenic Reiter, Andreas Gerhardus, and Jakob Runge. “Causal inference for temporal patterns”. arXiv preprint arXiv:2205.15149. 2022.
- [128] Yeonwoo Rho and Xiaofeng Shao. “Improving the bandwidth-free inference methods by prewhitening”. *Journal of Statistical Planning and Inference* 143.11 (2013), pp. 1912–1922.
- [129] Alessandro Rinaldo, Larry Wasserman, and Max G’Sell. “Bootstrapping and sample splitting for high-dimensional, assumption-lean inference”. *The Annals of Statistics* 47.6 (2019), pp. 3438–3469.
- [130] James Robins et al. “Higher order influence functions and minimax estimation of nonlinear functionals”. *Probability and Statistics: Essays in Honor of David A. Freedman* 2 (2008), pp. 335–422.
- [131] James Robins et al. “Semiparametric minimax rates”. *Electronic Journal of Statistics* 3 (2009), pp. 1305–1321.
- [132] James M. Robins et al. “Minimax estimation of a functional on a structured high-dimensional model”. *Annals of Statistics* 45.5 (2017), pp. 1951–1987.
- [133] Joseph P. Romano and Azeem M. Shaikh. “Inference for identifiable parameters in partially identified econometric models”. *Journal of Statistical Planning and Inference* 138.9 (2008), pp. 2786–2807.
- [134] Joseph P. Romano, Azeem M. Shaikh, and Michael Wolf. “A practical two-step method for testing moment inequalities”. *Econometrica: Journal of the Econometric Society* 82.5 (2014), pp. 1979–2002.
- [135] Murray Rosenblatt. “Independence and dependence”. *Proc. 4th Berkeley Symp. on Math. Statist. and Prob.* 2 (1961), pp. 431–443.
- [136] Reuven Y. Rubinstein and Dirk P. Kroese. *Simulation and the Monte Carlo method*. John Wiley and Sons, 2016.
- [137] Jakob Runge. “Causal network reconstruction from time series: from theoretical assumptions to practical estimation”. *Chaos: An Interdisciplinary Journal of Nonlinear Science* 28.7 (2018).
- [138] Jakob Runge. “Conditional independence testing based on a nearest-neighbor estimator of conditional mutual information”. *International Conference on Artificial Intelligence and Statistics* 84 (2018), pp. 938–947.
- [139] Jakob Runge et al. “Causal inference for time series”. *Nature Reviews Earth and Environment* 4.7 (2023), pp. 487–505.
- [140] Jakob Runge et al. “Detecting and quantifying causal associations in large nonlinear time series datasets”. *Science Advances* 5.11 (2019).
- [141] Jakob Runge et al. “Inferring causation from time series in Earth system sciences”. *Nature communications* 10.1 (2019).
- [142] Agathe Sadeghi, Achintya Gopal, and Mohammad Fesanghary. “Causal discovery from nonstationary time series”. *International Journal of Data Science and Analytics* 29.2 (2024), pp. 1–27.
- [143] Elena Saggioro et al. “Reconstructing regime-dependent causal relationships from observational time series”. *Chaos: An Interdisciplinary Journal of Nonlinear Science* 30.11 (2020).
- [144] Cyrill Scheidegger, Julia Hörrmann, and Peter Bühlmann. “The weighted generalised covariance measure”. *Journal of Machine Learning Research* 23.273 (2022), pp. 1–68.
- [145] Rajat Sen et al. “Model-powered conditional independence test”. *Advances in Neural Information Processing Systems* 30 (2017).
- [146] Sohan Seth and Jose C. Principe. “Assessing Granger non-causality using nonparametric measure of conditional independence”. *IEEE Transactions on Neural Networks and Learning Systems* 1 (2011), pp. 47–59.

- [147] Shalev Shaer, Gal Maman, and Yaniv Romano. “Model-X sequential testing for conditional independence via testing by betting”. *International Conference on Artificial Intelligence and Statistics* (2023), pp. 2054–2086.
- [148] Rajen D. Shah and Peter Bühlmann. “Goodness-of-fit tests for high dimensional linear models”. *Journal of the Royal Statistical Society Series B: Statistical Methodology* 80.1 (2018), pp. 113–135.
- [149] Rajen D. Shah and Jonas Peters. “The hardness of conditional independence testing and the generalised covariance measure”. *Annals of Statistics* 48.3 (2020), pp. 1514–1538.
- [150] Xiaofeng Shao. “A self-normalized approach to confidence interval construction in time series”. *Journal of the Royal Statistical Society Series B: Statistical Methodology* 72.3 (2010), pp. 343–366.
- [151] Xiaofeng Shao. “Self-normalization for time series: a review of recent developments”. *Journal of the American Statistical Association* 110.512 (2015), pp. 1797–1817.
- [152] Ali Shojaie and Emily B. Fox. “Granger causality: a review and recent advances”. *Annual Review of Statistics and Its Application* 9.1 (2022), pp. 289–319.
- [153] Peter Spirtes, Clark N. Glymour, and Richard Scheines. *Causation, Prediction, and Search*. MIT Press, 2001.
- [154] Ingo Steinwart and Marian Anghel. “Consistency of support vector machines for forecasting the evolution of an unknown ergodic dynamical system from observations with unknown noise”. *Annals of Statistics* 37.2 (2009), pp. 841–875.
- [155] Ingo Steinwart, Don Hush, and Clint Scovel. “Learning from dependent observations”. *Journal of Multivariate Analysis* 100.1 (2009), pp. 175–194.
- [156] Liangjun Su and Halbert White. “A consistent characteristic function-based test for conditional independence”. *Journal of Econometrics* 141.2 (2007), pp. 807–834.
- [157] Liangjun Su and Halbert White. “A nonparametric Hellinger metric test for conditional independence”. *Econometric Theory* 24.4 (2008), pp. 829–864.
- [158] Liangjun Su and Halbert White. “Testing conditional independence via empirical likelihood”. *Journal of Econometrics* 182.1 (2014), pp. 27–44.
- [159] Alex Tank et al. “Neural Granger causality”. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44.8 (2021), pp. 4267–4279.
- [160] Ryan J. Tibshirani et al. “Uniform asymptotic inference and the bootstrap after model selection”. *The Annals of Statistics* 46.3 (2018), pp. 1255–1287.
- [161] Michael Vogt. “Nonparametric regression for locally stationary time series”. *Annals of Statistics* 50.5 (2012), pp. 2601–2633.
- [162] Jonas Wahl, Urmi Ninad, and Jakob Runge. “Foundations of causal discovery on groups of variables”. *Journal of Causal Inference* 12.1 (2024).
- [163] Di Wang and Ruey S. Tsay. “Rate-optimal robust estimation of high-dimensional vector autoregressive models”. *The Annals of Statistics* 51.2 (2023), pp. 846–877.
- [164] Ian Waudby-Smith and Aaditya Ramdas. “Distribution-uniform anytime-valid inference”. arXiv preprint arXiv:2311.03343. 2023.
- [165] Ian Waudby-Smith et al. “Time-uniform central limit theory and asymptotic confidence sequences”. *Annals of Statistics* 52.6 (2024), pp. 2613–2640.
- [166] Michael Wieck-Sosa, Michel F. C. Haddad, and Aaditya Ramdas. “Identifying auxiliary indicators in unstable environments: are viral load distributions relevant for disease forecasting?” 2025.
- [167] Norbert Wiener. *Nonlinear problems in random theory*. MIT Press, 1966.
- [168] Kam Chung Wong, Zifan Li, and Ambuj Tewari. “Lasso guarantees for beta-mixing heavy-tailed time series”. *The Annals of Statistics* 48.2 (2020), pp. 1124–1142.
- [169] Wei Biao Wu. “Asymptotic theory for stationary processes”. *Statistics and its Interface* 4.2 (2011), pp. 207–226.

- [170] Wei Biao Wu. “Nonlinear system theory: another look at dependence”. *Proceedings of the National Academy of Sciences* 102.40 (2005), pp. 14150–14154.
- [171] Wei Biao Wu and Han Xiao. “Covariance matrix estimation in time series”. *Handbook of Statistics*. Vol. 30. Elsevier, 2012, pp. 187–209.
- [172] Wei-Biao Wu and Ying Nian Wu. “Performance bounds for parameter estimates of high-dimensional linear models with correlated errors”. *Electronic Journal of Statistics* 10.1 (2016), pp. 352–379.
- [173] Weichi Wu and Zhou Zhou. “Multiscale jump testing and estimation under complex temporal dynamics”. *Bernoulli* 30.3 (2024), pp. 2372–2398.
- [174] Jiaqi Xia, Yu Chen, and Xiao Guo. “Inference for high-dimensional linear models with locally stationary error processes”. *Journal of Time Series Analysis* 45.1 (2024), pp. 78–102.
- [175] Fang Xie and Zhijie Xiao. “Square-root LASSO for high-dimensional sparse linear systems with weakly dependent errors”. *Journal of Time Series Analysis* 39.2 (2018), pp. 212–238.
- [176] Fang Xie, Lihu Xu, and Youcai Yang. “Lasso for sparse linear regression with exponentially beta-mixing errors”. *Statistics and Probability Letters* 125 (2017), pp. 64–70.
- [177] Kashif Yousuf and Serena Ng. “Boosting high dimensional predictive regressions with time varying parameters”. *Journal of Econometrics* 224.1 (2021), pp. 60–87.
- [178] Bin Yu. “Rates of convergence for empirical processes of stationary mixing sequences”. *The Annals of Probability* (1994), pp. 94–116.
- [179] Danna Zhang and Wei Biao Wu. “Gaussian approximation for high dimensional time series”. *Annals of Statistics* 45.5 (2017), pp. 1895–1919.
- [180] Hao Zhang, Shuigeng Zhou, and Jihong Guan. “Causal discovery using regression-based conditional independence tests”. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence* (2017), pp. 1250–1256.
- [181] Hao Zhang et al. “Measuring conditional independence by independent residuals for causal discovery”. *ACM Transactions on Intelligent Systems and Technology* 10.5 (2019).
- [182] Kun Zhang et al. “Kernel-based conditional independence test and application in causal discovery”. *Conference on Uncertainty in Artificial Intelligence* 27 (2011), pp. 804–813.
- [183] Lu Zhang and Lucas Janson. “Floodgate: inference for model-free variable importance”. arXiv preprint arXiv:2007.01283. 2020.
- [184] Tianyu Zhang and Noah Simon. “A sieve stochastic gradient descent estimator for online non-parametric regression in Sobolev ellipsoids”. *The Annals of Statistics* 50.5 (2022), pp. 2848–2871.
- [185] Tianyu Zhang and Noah Simon. “Regression in tensor product spaces by the method of sieves”. *Electronic Journal of Statistics* 17.2 (2023), pp. 3660–3727.
- [186] Ting Zhang and Wei Biao Wu. “Time-varying nonlinear regression models: nonparametric estimation and model selection”. *Annals of Statistics* 43.2 (2015), pp. 741–768.
- [187] Xianyang Zhang and Guang Cheng. “Bootstrapping high dimensional time series”. arXiv preprint arXiv:1406.1037. 2014.
- [188] Yi Zhang and Xiaofeng Shao. “Another look at bandwidth-free inference: a sample splitting approach”. *Journal of the Royal Statistical Society Series B: Statistical Methodology* 86.1 (2024), pp. 246–272.
- [189] Zhou Zhou. “Heteroscedasticity and autocorrelation robust structural change detection”. *Journal of the American Statistical Association* 108.502 (2013), pp. 726–740.
- [190] Zhou Zhou and Wei Biao Wu. “Local linear quantile estimation for nonstationary time series”. *The Annals of Statistics* 37.5B (2009), pp. 2696–2729.

A Extensions

A.1 Alternative test statistics

Consider the test statistic

$$S_{n,p}^*(\hat{\mathbf{R}}_n) = \left\| \frac{1}{\sqrt{T_{n,L}}} \sum_{t \in \mathcal{T}_{n,L}} \hat{\mathbf{R}}_{t,n} \right\|_p,$$

based on the ℓ_p -norm ($p \geq 2$) of the full sum of residual products. For example, we can use the test statistics

$$S_{n,\infty}^*(\hat{\mathbf{R}}_n) = \left\| \frac{1}{\sqrt{T_{n,L}}} \sum_{t \in \mathcal{T}_{n,L}} \hat{\mathbf{R}}_{t,n} \right\|_\infty, \quad S_{n,2}^*(\hat{\mathbf{R}}_n) = \left\| \frac{1}{\sqrt{T_{n,L}}} \sum_{t \in \mathcal{T}_{n,L}} \hat{\mathbf{R}}_{t,n} \right\|_2.$$

Crucially, the full sum test statistic $S_{n,p}^*(\hat{\mathbf{R}}_n)$ will not have power against alternatives in which the time-averages of the time-varying expected conditional covariances are close to zero (e.g. positive during the first half of times, and negative during the second half). On the other hand, the maximum partial sum test statistic $S_{n,p}(\hat{\mathbf{R}}_n)$ from (10) *does* have power against these alternatives. If the time-varying expected conditional covariances are suspected to consistently maintain the same sign (whether positive or negative), then users might be able to gain *some* power by using $S_{n,p}^*(\hat{\mathbf{R}}_n)$, although we emphasize that $S_{n,p}(\hat{\mathbf{R}}_n)$ will also have power against these alternatives. However, in settings where we have little prior knowledge about the time-varying expected conditional covariances between the nonstationary processes under alternatives, then the maximum partial sum test statistic $S_{n,p}(\hat{\mathbf{R}}_n)$ should be used because it has power against a wider range of alternatives. For similar reasons, we recommend using $S_{n,p}(\hat{\mathbf{R}}_n)$ when conducting automated multiple conditional independence testing (e.g. for screening out irrelevant time series in a large database of possible forecasting signals).

It is perhaps most intuitive to frame the problem in the following way. Consider the time-varying partially linear model

$$\mathbb{E}_P(Y_{t,n,j,b} | X_{t,n,i,a}, \mathbf{Z}_{t,n}) = \beta_{P,t,n,m} X_{t,n,i,a} + h_{P,t,n,j,b}(\mathbf{Z}_{t,n}),$$

for some function $h_{P,t,n,j,b}(\cdot)$. When the time-varying conditional expectation $\mathbb{E}_P(Y_{t,n,j,b} | X_{t,n,i,a}, \mathbf{Z}_{t,n})$ is assumed to have this time-varying partially linear form, the time-varying coefficient $\beta_{P,t,n,m}$ is equal to the expected conditional covariance of $X_{t,n,i,a}$ and $Y_{t,n,j,b}$ given $\mathbf{Z}_{t,n}$ divided by the expected conditional variance of $X_{t,n,i,a}$ given $\mathbf{Z}_{t,n}$; see Robins et al. [131] and Hines et al. [71] for more discussion. If domain knowledge suggests that the time-varying coefficients $(\beta_{P,t,n,m})_{m=(i,j,a,b) \in \mathcal{D}_n}$ consistently maintain the same sign over time $t \in \mathcal{T}_n$, then the full sum test statistic $S_{n,p}^*(\hat{\mathbf{R}}_n)$ can be used to gain some power. Otherwise, if we cannot make this assumption, then use the maximum partial sum test statistic $S_{n,p}(\hat{\mathbf{R}}_n)$ which has power against a broader spectrum of alternatives.

A.2 Cyclostationary processes

The general triangular array framework from Section 3 also allows for nonstationary processes that exhibit some form of repetition over time, such as periodic stationarity or cyclostationarity [17, 117, 18, 58, 59, 112]. We emphasize that these types of nonstationary processes are not necessarily locally stationary. The theoretical justification of the dGCM test from Theorem 3.1 requires that we improve our estimates of the time-varying regression functions as n grows. See Remark 2.1 in Chen, Smetanina, and Wu [27] and the preceding discussion about time-varying regression with periodic stationary or cyclostationary processes. Also, see Subsection 2.5.1 of Bonnerjee, Karmakar, and Wu [20] for a discussion of how strong Gaussian approximations for nonstationary nonlinear processes with causal representations as in Subsections 3.1, 3.3, 3.4 can be used with cyclostationary processes. Ideally, we would like to be able to handle even more complex forms of nonstationarity than cyclostationarity. See Gardner, Napolitano, and Paura [59] and Napolitano [111] for generalizations of this concept.

A.3 Simulation-and-regression for nonstationary processes

The general triangular array framework from Section 3 can also be used with simulation-and-regression approaches for estimating the time-varying regression functions. Suppose we have access to a black-box simulator which can be used to generate realistic paths of (X, Z) . Naturally, this simulation-based

approach assumes that we either know the parameters of the simulator, or how to estimate them using an appropriate technique [34, 107, 61, 56]. The main idea of this approach is to simulate s paths of (X, Z) , then fit separate regression models, such as XGBoost [28], LightGBM [83], or random forests [22], for each time t by using the observations across the s iid simulated paths. To obtain the residuals for the time-varying regression of X on Z , the fitted regression models for each time t can be used with the observed realization of (X, Z) . The residuals for the time-varying regression of Y on Z can be obtained as in Section 4 or Subsection A.2.

The asymptotic arguments can be based on letting the number of simulations s grow with n , where n can be linked to the sample size (e.g. sampling frequency and/or duration of time) and number of dimensions. We can also allow n to be linked to the quality of the simulations, so that as n grows we can generate more realistic simulations — perhaps at a higher computational cost — and the simulator can be seen as converging in some sense to the true data generating process. Note that if the simulator can generate paths for (X, Y, Z) , then conditional independence tests which require multiple realizations of a nonstationary process, such as [105, 100, 98], can be used. In contrast, our proposed approach only requires a simulator for (X, Z) and a single realization of Y .

The main advantage of this simulation-based approach is that it leverages domain knowledge about (X, Z) to obtain better estimates of the time-varying regression functions without assuming anything about Y . For example, stochastic simulators can be used to generate paths of climatic variables, such as precipitation, surface temperature, surface water vapor, and ozone. Using the approach described above, we can identify conditional dependencies between these simulatable climatic variables (X, Z) and another process Y that is harder to model (e.g. commodity prices, consumer behavior, flu cases).

In this setting, it may also be possible to develop a conditional independence test using a model-X approach [24] based on simulation-based conditional density estimation. Nevertheless, the dGCM test can still be a very practical choice for this setting, particularly when it is much more feasible to do simulation-and-regression than simulation-based conditional density estimation. For instance, when the simulator’s computational demands make it impossible to generate a large number of paths for a high-dimensional process (X, Z) .

A.4 Simplifications under stationarity

Throughout this paper, we have completely avoided the assumption of stationarity. However, it is worth explaining how things would simplify if we are willing to assume that the processes are stationary. Overall, the takeaway is that the original GCM test from Shah and Peters [149] would require minimal modifications.

To begin, suppose we have $n \in \mathbb{N}$ observations of a stationary mixing time series, so that the regression functions are time-invariant. Further, suppose that the errors are iid. The statistical guarantees of many machine learning algorithms and statistical models, such as support vector machines [154, 155, 67], random forests [60, 42], lasso [168], and high-dimensional vector autoregressive models [163], have been studied in the context of stationary mixing time series with iid errors. Over the last decade, the literature on statistical learning theory for time series has been able to move beyond the restrictive assumptions of stationarity and mixing [178, 168, 80, 2] (or asymptotic stationarity [1]) by describing nonstationarity in terms of discrepancy measures [91, 92, 90, 69, 110]. This literature has recently considered new notions of learnability for general non-iid stochastic processes [43] and conditions under which learning from general non-iid stochastic processes is possible [68].

Zhang and Wu [186] considered the setting in which the regression functions are time-varying and the errors are iid. Since the errors are iid, a multiplier bootstrap testing procedure can be justified by the Gaussian approximation from Chernozhukov, Chetverikov, and Kato [30] which was used by Shah and Peters [149]. Hence, the resulting test would be very similar to the original GCM test for the iid setting from Shah and Peters [149]. The main difference is that there can be time-lagged conditional dependencies in the stationary time series setting.

Suppose that the observed processes are temporally dependent (e.g. some form of mixing) and stationary so that the regression functions are time-invariant as before, but the errors are also temporally dependent. The guarantees of the lasso and vector autoregressive models are fairly well-studied in this setting. Basu and Michailidis [10] investigated high-dimensional vector autoregressive models with serially correlated errors. Gupta [66] and Xie and Xiao [175] studied the lasso with errors satisfying various weak dependence conditions. Peng, Zhu, and Zhong [121] and Xie, Xu, and Yang [176] studied the lasso with ϕ -mixing and β -mixing errors, respectively. Wu and Wu [172] studied the guarantees of

the lasso in the setting with temporally dependent errors by using the functional dependence measure of Wu [170].

In the serially correlated error setting, the key difference with the GCM test from Shah and Peters [149] is that one must use a suitable Gaussian approximation result to justify a multiplier bootstrap-type testing procedure. See Chang, Chen, and Wu [26] for a comprehensive overview of Gaussian approximations for dependent data. Chernozhukov, Chetverikov, and Kato [32, 31] investigated a block multiplier bootstrap under a β -mixing assumption, and Zhang and Cheng [187] explored a wild multiplier bootstrap under the functional dependence measure of Wu [170]. Also, Zhang and Wu [179] discuss estimators for the long-run covariance matrix so that their Gaussian approximation for high-dimensional time series can be applied in practice. See Wu and Xiao [171] and Wu [169] for more discussion about long-run covariance matrix estimation for stationary time series.

One could also have time-invariant regression functions with errors that are nonstationary and temporally dependent. For instance, Xia, Chen, and Guo [174] studies the lasso with locally stationary errors. However, the statistical guarantees of other machine learning algorithms and statistical models have not been studied in this setting. If the process of error products is mean-nonstationary (i.e. time-varying expected conditional covariance) under alternatives, then the same test statistics from Subsection 2.4 can be used. Otherwise, if domain knowledge suggests that the time-varying expected conditional covariances usually maintain the same sign, then the test statistics from Subsection A.1 can be used.

To recap, in this subsection we considered how the assumption of stationarity would vastly simplify the problem. We find that the original GCM test for the iid setting from Shah and Peters [149] can be adapted to the stationary time series setting by making the previously mentioned changes. In contrast, we consider the much more complicated setting in which the observed processes can be nonstationary and temporally dependent, the regression functions can vary over time, and the error processes can be nonstationary and temporally dependent. We emphasize that our dGCM test can be used with stationary processes and iid sequences, which are special cases of the general framework from Section 3.

A.5 Additional tests for locally stationary processes

In this subsection, we discuss three conditional independence tests for locally stationary processes that we did not pursue in this paper. Crucially, the test statistics below require local long-run covariance estimation. Most local long-run covariance estimators use kernel smoothing, and therefore require selecting bandwidths. Unfortunately, test statistics that use kernel smoothing can be very sensitive to the choice of the bandwidths, which can be hard to select in practice. Inspired by the success of bandwidth-free approaches in other areas of time series analysis [99, 150, 128, 151, 188], we designed the dGCM test so that it does not require local long-run covariance estimation and therefore avoids kernel smoothing.

Recall the notation for the locally stationary setting from Section 4. To begin, let us translate the “weak” conditional independence criterion of Daudin [41] into the locally stationary setting as follows. For some $n \in \mathbb{N}$, $u \in \mathcal{U}_n$, $(i, j, a, b) \in \mathcal{D}_n$, $t \in \mathcal{T}_n$, if

$$\tilde{X}_{\lfloor un \rfloor, n, i, a}(u) \perp\!\!\!\perp \tilde{Y}_{\lfloor un \rfloor, n, j, b}(u) \mid \tilde{Z}_{\lfloor un \rfloor, n}(u),$$

then

$$\mathbb{E}_P[\phi(\tilde{X}_{\lfloor un \rfloor, n, i, a}(u), \tilde{Z}_{\lfloor un \rfloor, n}(u))\varphi(\tilde{Y}_{\lfloor un \rfloor, n, j, b}(u), \tilde{Z}_{\lfloor un \rfloor, n}(u))] = 0,$$

for all functions

$$\phi \in L^2_{\tilde{X}_{\lfloor un \rfloor, n, i, a}(u), \tilde{Z}_{\lfloor un \rfloor, n}(u)}, \quad \varphi \in L^2_{\tilde{Y}_{\lfloor un \rfloor, n, j, b}(u), \tilde{Z}_{\lfloor un \rfloor, n}(u)},$$

such that

$$\begin{aligned} \mathbb{E}_P[\phi(\tilde{X}_{\lfloor un \rfloor, n, i, a}(u), \tilde{Z}_{\lfloor un \rfloor, n}(u)) \mid \tilde{Z}_{\lfloor un \rfloor, n}(u)] &= 0, \\ \mathbb{E}_P[\varphi(\tilde{Y}_{\lfloor un \rfloor, n, j, b}(u), \tilde{Z}_{\lfloor un \rfloor, n}(u)) \mid \tilde{Z}_{\lfloor un \rfloor, n}(u)] &= 0. \end{aligned}$$

Hence, the corresponding *local* expected conditional covariance

$$\rho_{P, n, m}(u) = \mathbb{E}_P[\text{Cov}_P(\tilde{X}_{\lfloor un \rfloor, n, i, a}(u), \tilde{Y}_{\lfloor un \rfloor, n, j, b}(u) \mid \tilde{Z}_{\lfloor un \rfloor, n}(u))],$$

is equal to zero for $m = (i, j, a, b) \in \mathcal{D}_n$.

First, consider the global null hypothesis of conditional independence

$$\tilde{X}_{\lfloor un \rfloor, n, i, a}(u) \perp\!\!\!\perp \tilde{Y}_{\lfloor un \rfloor, n, j, b}(u) \mid \tilde{Z}_{\lfloor un \rfloor, n}(u) \text{ for all } u \in \mathcal{U}_n, \text{ for all } (i, j, a, b) \in \mathcal{D}_n. \quad (15)$$

In the univariate setting, \mathcal{D}_n simply consists of one dimension/time-offset tuple as in Subsection 2.2. Also, note that this hypothesis can be extended to the group of time series setting as discussed in Subsection 2.2. Note that the null hypothesis (15) implies the null hypothesis (14), so the process of error products from the time-varying nonlinear regressions of $(X_{t,n,i,a})_{t \in \mathcal{T}_n}$ on $(\mathbf{Z}_{t,n})_{t \in \mathcal{T}_n}$ and $(Y_{t,n,j,b})_{t \in \mathcal{T}_n}$ on $(\mathbf{Z}_{t,n})_{t \in \mathcal{T}_n}$ will still have mean zero as in Section 4. To test for the null hypothesis (15), we could, for example, use the test statistic

$$\sup_{u \in \mathcal{U}_n} \left\| \frac{1}{\sqrt{T_n}} \sum_{t=\mathbb{T}_n^-}^{\lfloor un \rfloor} (\hat{\Sigma}_{t,n}^R)^{-1/2} \hat{\mathbf{R}}_{t,n} \right\|_p,$$

based on some ℓ_p norm ($p \geq 2$) of the studentized partial sum process, where $\hat{\Sigma}_{t,n}^R$ is an estimate of the local long-run covariance matrix at time t . The theoretical guarantees for the test based on this test statistic would utilize the recent results from Mies [108] about strong Gaussian approximations with random multipliers.

Second, it is possible to develop a test for the *local* null hypothesis of conditional independence

$$\tilde{X}_{\lfloor un \rfloor, n, i, a}(u) \perp\!\!\!\perp \tilde{Y}_{\lfloor un \rfloor, n, j, b}(u) \mid \tilde{Z}_{\lfloor un \rfloor, n}(u) \text{ for all } (i, j, a, b) \in \mathcal{D}_n, \quad (16)$$

for a *particular* rescaled time $u \in \mathcal{U}_n$ (i.e. instead of for all \mathcal{U}_n as in (15)) by using, for example, the test statistic

$$\max_{m=(i,j,a,b) \in \mathcal{D}_n} \left| \frac{1}{\sqrt{T_n h_{n,m}}} \sum_{t \in \mathcal{T}_n} K\left(\frac{t/n - u}{h_n}\right) \hat{R}_{t,n,m} \right| / \hat{\sigma}_{n,m}^R(u),$$

for some bandwidths $h_{n,m} \rightarrow 0$ and local long-run variance estimates $(\hat{\sigma}_{n,m}^R(u))^2$. The main idea for this local conditional independence test is that since $\mathbb{E}_P(\tilde{R}_{P,\lfloor un \rfloor, n, m}(u)) = 0$ under the null and the process of error products is “approximately stationary” over short periods of time, we can expect that the means of $R_{P,t,n,m} = \tilde{R}_{P,t,n,m}(t/n)$ for rescaled times t/n near u are also close to zero under the null. We can make this mathematically precise by using the technical tools developed for nonlinear locally stationary processes from Dahlhaus, Richter, and Wu [40].

Third, it is also possible to simultaneously test whether conditional independence

$$\tilde{X}_{\lfloor un \rfloor, n, i, a}(u) \perp\!\!\!\perp \tilde{Y}_{\lfloor un \rfloor, n, j, b}(u) \mid \tilde{Z}_{\lfloor un \rfloor, n}(u) \text{ for all } (i, j, a, b) \in \mathcal{D}_n, \quad (17)$$

holds at *each* rescaled time $u \in \mathcal{U}_n$ (i.e. instead of for a particular $u \in \mathcal{U}_n$ as in (16)). This can be done by creating *simultaneous confidence bands* (i.e. over time) for expected conditional covariance curves $(\rho_{P,n,m}(u))_{u \in \mathcal{U}_n}$ for each $m = (i, j, a, b) \in \mathcal{D}_n$. Depending on whether or not estimates of the local long-run variances $(\hat{\sigma}_{n,m}^R(u))^2$ are used, these simultaneous confidence bands will have time-varying or time-invariant widths, respectively. The main idea is that the local null hypothesis of conditional independence at rescaled time $u \in \mathcal{U}_n$ for some dimension/time-offset tuple $m = (i, j, a, b) \in \mathcal{D}_n$ can be rejected if zero is not included in the corresponding confidence interval for the local expected conditional covariance $\rho_{P,n,m}(u)$. This can be done using similar arguments as Bai and Wu [7], which focuses on inferring time-varying correlation curves. However, due to the problem of post-selection inference [87], this would require either stronger assumptions (e.g. Donsker-type), data decomposition techniques (e.g. splitting, fission, or thinning) for nonstationary time series, or two independent realizations of the same nonstationary process — rarely possible outside of experimental settings.

An approach for inferring expected conditional covariance curves would have a range of applications outside of testing for conditional independence, since this functional frequently appears in the causal inference literature [84, 130, 131, 95, 132, 114]. We suspect that similar approaches can be used to infer curves based on other functionals of interest in causal inference. Hence, this line of work would be of significant interest to the emerging field of time series causal inference [143, 127, 139, 140, 137]. Lastly, we note that it may be possible to extend the tests discussed in this subsection to the piecewise locally stationary setting (see Subsection A.6), however we leave the details for future work.

A.6 Piecewise locally stationary processes

We briefly describe how to extend the Sieve-dGCM test from Section 4 from locally stationary processes [36, 190, 37, 40] to a more general class of nonstationary processes known as *piecewise locally stationary* (PLS) processes introduced in Zhou [189]. Specifically, the class of PLS processes generalizes the stochastic Lipschitz nonstationarity condition from Assumption 4.6 by allowing for finitely many breakpoints [189, 173, 44]. We emphasize that PLS processes are included in the even more general class of nonstationary processes from Section 3 with the total variation-type nonstationarity condition from Assumption 3.6.

The main idea is to identify the breakpoints, fit a separate sieve model on each locally stationary segment, and run Algorithm 1 on *all* the residuals. If the breakpoints are known exactly, then the same arguments can be used to show that the sieve time-varying regression estimators achieve the required convergence rates (i.e. within each locally stationary segment). If the breakpoints must be identified, then our arguments must be extended to account for this. As far as we know, Wu and Zhou [173] is the most relevant work on identifying breakpoints for PLS processes. We leave the full details of this extension for future work.

A.7 Weakening the assumptions on the error processes

In Assumption 3.5, we assume that there are distribution-uniform upper bounds on the L^∞ norms and L^∞ functional dependence measures of the error processes. We use this assumption to show inequality (23) in the proof of Theorem 3.1. Afterwards, we use the time-uniform convergence rates for the time-varying regression estimators to show Step 1.2 in the proof of Theorem 3.1. It is possible to weaken the assumptions imposed on the error processes by making stronger assumptions about the time-varying regression estimators, or more complicated assumptions about the terms in (23). Instead, we opt for simpler assumptions on the errors and prediction errors for the sake of transparency.

Lastly, we note that the Sieve-dGCM test from Section 4 performs well even when the error processes violate Assumption 3.5, at least in the settings we considered for our simulations in Section 5. To satisfy Assumption 3.5, we can make minor modifications to the data generating processes in Section 5. For example, by replacing the Gaussian error processes with, say, truncated Gaussian error processes.

B Distribution-Uniform Theory

In this section, we state distribution-uniform versions of the results from Mies and Steland [109]. All of the results in this section can be applied to general triangular array frameworks for high-dimensional nonstationary nonlinear processes, such as locally stationary processes.

B.1 Literature review of distribution-uniform inference

We briefly review prior work on distribution-uniform inference. First, we discuss the conditional independence testing literature. Recently, there has been a great deal of work on distribution-uniform conditional independence testing frameworks due to the hardness result and conditional independence testing framework from Shah and Peters [149]. For instance, Lundborg, Shah, and Peters [100] introduced many distribution-uniform convergence results for separable Banach and Hilbert spaces. Recently, Christgau, Petersen, and Hansen [33] introduced a distribution-uniform “conditional local independence” testing framework for the setting in which n realizations of a point process are observed. Christgau, Petersen, and Hansen [33] also introduce a distribution-uniform version of Rebolledo’s martingale central limit theorem [126] and extend many distribution-uniform convergence results from Lundborg, Shah, and Peters [100] to metric spaces.

Second, we mention some relevant work from the literature on anytime-valid inference. Recently, Waudby-Smith and Ramdas [164] introduced a distribution-uniform strong (almost-sure) Gaussian approximation for the full sum of iid random variables. The work in Waudby-Smith and Ramdas [164] is motivated by prior work on asymptotic anytime-valid inference from Waudby-Smith et al. [165], in which the authors defined the concept of an “asymptotic confidence sequence”. In particular, Waudby-Smith et al. [165] introduced asymptotic confidence sequences for iid random variables and a Lindeberg-type asymptotic confidence sequence which can capture time-varying means under martingale dependence.

Third, we discuss other areas in which distribution-uniform inference is studied under different names. There is a vast literature discussing the importance of distribution-uniform inference under the name of “honest” or “uniform” inference, see [94, 82, 160, 129, 88]. Also, there is a plethora of literature on distribution-uniform moment inequality testing [77, 133, 5, 6, 4, 134]. Most recently, Li, Liao, and Zhou [93] developed a distribution-uniform test for general functional inequalities which admits conditional moment inequalities as a special case. In the supplemental appendix, Li, Liao, and Zhou [93] introduce a distribution-uniform strong Gaussian approximation for the full sum of a high-dimensional mixingale.

B.2 Distribution-uniform strong Gaussian approximation

To begin, let us introduce the setting rigorously. Let Ω be a sample space, \mathcal{B} the Borel sigma-algebra, and (Ω, \mathcal{B}) a measurable space. For fixed $n \in \mathbb{N}$, let (Ω, \mathcal{B}) be equipped with a family of probability measures $(\mathbb{P}_P)_{P \in \mathcal{P}_n}$ so that the distribution of the high-dimensional stochastic system

$$(G_{t,n}(\mathcal{H}_s))_{t \in [n], s \in \mathbb{Z}},$$

or, in the locally stationary setting

$$(\tilde{G}_n(u, \mathcal{H}_s))_{u \in [0,1], s \in \mathbb{Z}},$$

under \mathbb{P}_P is $P \in \mathcal{P}_n$. Here $\mathcal{H}_t = (\eta_t, \eta_{t-1}, \dots)$, where $(\eta_t)_{t \in \mathbb{Z}}$ is a sequence of iid random vectors with dimension $d^\eta = d_n^\eta$, and $G_{t,n} : (\mathbb{R}^{d_n^\eta})^\infty \rightarrow \mathbb{R}^{d_n}$ is a measurable function — where we endow $(\mathbb{R}^{d_n^\eta})^\infty$ with the σ -algebra generated by all finite projections. For each $t \in [n]$, $G_{t,n}(\mathcal{H}_s)$ is a well-defined high-dimensional random vector for every $s \in \mathbb{Z}$, and $(G_{t,n}(\mathcal{H}_s))_{s \in \mathbb{Z}}$ is a high-dimensional stationary ergodic process.

For each $n \in \mathbb{N}$, write the \mathbb{R}^{d_n} -valued process of interest as $(W_{t,n})_{t \in [n]}$. We assume that for each $n \in \mathbb{N}$ and $t \in [n]$, the random vector $W_{t,n}$ has a causal representation; that is, it can be represented as a measurable function of these iid random vectors

$$W_{t,n} = G_{t,n}(\mathcal{H}_t).$$

Similarly, for the causal representations in the locally stationary setting, the measurable function $\tilde{G}_n(u, \cdot) : (\mathbb{R}^{d_n^\eta})^\infty \rightarrow \mathbb{R}^{d_n}$ is defined for each rescaled time $u \in [0, 1]$, and we assume that

$$W_{t,n} = \tilde{G}_n(t/n, \mathcal{H}_t).$$

We can use the results in this section for locally stationary processes by writing

$$G_{t,n}(\mathcal{H}_t) = \tilde{G}_n(t/n, \mathcal{H}_t).$$

The family of probability measures $(\mathbb{P}_P)_{P \in \mathcal{P}_n}$ is defined with respect to the same measurable space (Ω, \mathcal{B}) , but need not have the same dominating measure. Denote a family of probability spaces by $(\Omega, \mathcal{B}, \mathbb{P}_P)_{P \in \mathcal{P}_n}$. When we say that the process $(W_{t,n})_{t \in [n]}$ is defined on the collection of probability spaces $(\Omega, \mathcal{B}, \mathbb{P}_P)_{P \in \mathcal{P}_n}$ for some $n \in \mathbb{N}$, we mean that $(W_{t,n})_{t \in [n]}$ is defined on the probability space $(\Omega, \mathcal{B}, \mathbb{P}_P)$ for each $P \in \mathcal{P}_n$.

Note that the causal representations in this paper use sequences of iid random vectors, whereas the causal representations in Mies and Steland [109] use sequences of iid $\text{Unif}[0, 1]$ random variables. The same arguments used in Mies and Steland [109] can be applied when using our formulation of the causal representations with iid random vectors. The only reason we write the causal representations in this way is for the sake of clarity.

In fact, standard results in probability theory imply that the causal representations based on measurable functions of sequences of iid $\text{Unif}[0, 1]$ random variables as in Mies and Steland [109] are already sufficiently general. For example, see Kallenberg [79] Lemma 4.21, Lemma 4.22, and the surrounding discussion. More specifically, the causal representations with sequences of iid $\text{Unif}[0, 1]$ random variables can express the causal representations with sequences of random vectors by including compositions with additional measurable functions for (1) replicating each of the iid $\text{Unif}[0, 1]$ random variables, and (2) inverse sampling via products of conditional distributions; see Section 2.5 of Rubinstein and Kroese [136] on random vector generation.

Next, we define our measure of temporal dependence for the process. For the following definition, let $(\tilde{\eta}_t)_{t \in \mathbb{Z}}$ be an iid copy of $(\eta_t)_{t \in \mathbb{Z}}$ and denote

$$\tilde{\mathcal{H}}_{t,j} = (\eta_t, \dots, \eta_{j+1}, \tilde{\eta}_{t-j}, \eta_{t-j-1}, \dots),$$

to be \mathcal{H}_t with the j -th input in the past η_{t-j} replaced with the iid copy $\tilde{\eta}_{t-j}$.

Definition B.1 (Functional dependence measure). *For $n \in \mathbb{N}$, $P \in \mathcal{P}_n$, $t \in \mathcal{T}_n$, define the functional dependence measure of $W_{t,n} = G_{t,n}(\mathcal{H}_t)$ as*

$$\theta_{P,t,n}(j, q, r) = (\mathbb{E}_P \|G_{t,n}(\mathcal{H}_t) - G_{t,n}(\tilde{\mathcal{H}}_{t,j})\|_r^q)^{\frac{1}{q}},$$

with $h \in \mathbb{N}_0$, $q \geq 1$, $r \geq 1$.

We state the following distribution-uniform assumptions about the temporal dependence and non-stationarity of the process for some collections of distributions \mathcal{P}_n for some $n \in \mathbb{N}$. Note that we write the time of the input sequence as 0 when it does not matter due to stationarity.

Assumption B.1 (Distribution-uniform decay of temporal dependence). *We assume that there exist $\beta > 0$, $q \geq 2$ and a constant $\Theta_n > 0$, such that for all times $t \in [n]$ it holds*

$$\sup_{P \in \mathcal{P}_n} \theta_{P,t,n}(j, q, r) \leq \Theta_n \cdot (j \vee 1)^{-\beta},$$

for $j \geq 0$, and that

$$\sup_{P \in \mathcal{P}_n} (\mathbb{E}_P \|G_{t,n}(\mathcal{H}_0)\|_2^q)^{1/q} \leq \Theta_n.$$

Assumption B.2 (Distribution-uniform total variation condition for nonstationarity). *Recall Θ_n from Assumption B.1. Assume that there exists some $\Gamma_n \geq 1$ such that*

$$\sup_{P \in \mathcal{P}_n} \left(\sum_{t=2}^n (\mathbb{E}_P \|G_{t,n}(\mathcal{H}_0) - G_{t-1,n}(\mathcal{H}_0)\|_2^2)^{\frac{1}{2}} \right) \leq \Gamma_n \cdot \Theta_n.$$

Note that the assumptions regarding the temporal dependence and nonstationarity of the process of error products, as stated in Subsection 3.4, ensure that both Assumptions B.1 and B.2 hold for each $n \in \mathbb{N}$. Furthermore, since the assumptions in Subsection 4.5 are strictly stronger than those in Subsection 3.4, the results from this section can be applied to the process of error products in both Sections 3 and 4.

Define the two rates

$$\chi(q, \beta) = \begin{cases} \frac{q-2}{6q-4}, & \beta \geq \frac{3}{2}, \\ \frac{(\beta-1)(q-2)}{q(4\beta-3)-2}, & \beta \in (1, \frac{3}{2}), \end{cases}$$

and

$$\xi(q, \beta) = \begin{cases} \frac{q-2}{6q-4}, & \beta \geq 3, \\ \frac{(\beta-2)(q-2)}{(4\beta-6)q-4}, & \frac{3+\frac{2}{q}}{1+\frac{2}{q}} < \beta < 3, \\ \frac{1}{2} - \frac{1}{\beta}, & 2 < \beta \leq \frac{3+\frac{2}{q}}{1+\frac{2}{q}}, \end{cases}$$

which will appear in the results in this section. In general, the Gaussian approximation allows the dimensions to grow as $d_n = O(n^{\frac{1-\delta}{1+\frac{1}{2\xi(q,\beta)}}})$ for some $\delta > 0$. In the limiting case when $\beta \geq 3$ and $q \rightarrow \infty$, this corresponds to $d_n = O(n^{\frac{1}{4}-\delta'})$ for some $\delta' > 0$.

Let us briefly recall some notation used in the main text. Recall $\bar{\beta}^R > 3$, $\bar{q}^R > 4$ from Assumption 3.5, as well as the number of times T_n and the number of dimension/time-offset combinations D_n from Subsection 2.1. For the dGCM test, we allow $D_n = O(T_n^{r(\bar{q}^R, \bar{\beta}^R)})$ where

$$r(\bar{q}^R, \bar{\beta}^R) = \min \left(\frac{1-\delta}{1+\frac{1}{2\xi(\bar{q}^R, \bar{\beta}^R)}}, \frac{1}{6} \right), \quad (18)$$

for some $\delta > 0$. The limiting factor that leads to this requirement is *not* from the strong Gaussian approximation, but rather to ensure that the convergence rate requirements can be achieved by the time-varying nonparametric regression estimators.

The following result is a distribution-uniform version of the strong Gaussian approximation from Theorem 3.1 in Mies and Steland [109].

Lemma B.1. For some sample size $n \in \mathbb{N}$ and collection of distributions \mathcal{P}_n for the stochastic system $(G_{t,n}(\mathcal{H}_s))_{t \in [n], s \in \mathbb{Z}}$, suppose that Assumption B.1 is satisfied for \mathcal{P}_n with some $q > 2$, $\beta > 1$ and constant $\Theta_n > 0$. Let the \mathbb{R}^{d_n} -valued process $(W_{t,n})_{t \in [n]}$ be defined on the collection of probability spaces $(\Omega, \mathcal{B}, \mathbb{P}_P)_{P \in \mathcal{P}_n}$ so that $W_{t,n} = G_{t,n}(\mathcal{H}_t)$ with $\mathbb{E}_P(W_{t,n}) = 0$ for each time $t \in [n]$ and distribution $P \in \mathcal{P}_n$. Also, suppose the dimension $d_n < cn$ for some constant $c > 0$. Then, on a potentially enriched collection of probability spaces $(\Omega', \mathcal{B}', \mathbb{P}'_P)_{P \in \mathcal{P}_n}$, there exist random vectors $(W'_{t,n})_{t \in [n]}$ with the same distribution as $(W_{t,n})_{t \in [n]}$ for each $P \in \mathcal{P}_n$, and independent Gaussian random vectors $(V'_{t,n})_{t \in [n]}$ with $\mathbb{E}_P(V'_{t,n}) = 0$ for each $t \in [n]$, $P \in \mathcal{P}_n$, such that

$$\sup_{P \in \mathcal{P}_n} \left(\mathbb{E}_P \max_{k \leq n} \left\| \frac{1}{\sqrt{n}} \sum_{t=1}^k (W'_{t,n} - V'_{t,n}) \right\|_2^2 \right)^{\frac{1}{2}} \leq K \Theta_n \sqrt{\log(n)} \left(\frac{d_n}{n} \right)^{\chi(q,\beta)}$$

for some universal constant K depending only on q , c , and β .

If $\beta > 2$, then the local long-run covariance matrix $\Sigma_{P,t,n} = \sum_{h=-\infty}^{\infty} \text{Cov}_P(G_{t,n}(\mathcal{H}_0), G_{t,n}(\mathcal{H}_h))$ is well-defined for each $t \in [n]$, $P \in \mathcal{P}_n$ by Lemma B.5. If Assumption B.2 is also satisfied for \mathcal{P}_n , then on $(\Omega', \mathcal{B}', \mathbb{P}'_P)_{P \in \mathcal{P}_n}$ there exist random vectors $(W'_{t,n})_{t \in [n]}$ which have the same distribution as $(W_{t,n})_{t \in [n]}$ for each $P \in \mathcal{P}_n$, and independent Gaussian random vectors $(V^*_{t,n})_{t \in [n]}$ where $V^*_{t,n} \sim \mathcal{N}(0, \Sigma_{P,t,n})$ for each $t \in [n]$, $P \in \mathcal{P}_n$, such that

$$\sup_{P \in \mathcal{P}_n} \left(\mathbb{E}_P \max_{k \leq n} \left\| \frac{1}{\sqrt{n}} \sum_{t=1}^k (W'_{t,n} - V^*_{t,n}) \right\|_2^2 \right)^{\frac{1}{2}} \leq K \Theta_n \Gamma_n^{\frac{1}{2} \frac{\beta-2}{\beta-1}} \sqrt{\log(n)} \left(\frac{d_n}{n} \right)^{\xi(q,\beta)}$$

for some universal constant K depending only on q , c , and β .

Proof of Lemma B.1: Assumptions B.1 and B.2 are distribution-uniform versions of conditions (G.1) and (G.2) from Mies and Steland [109]. Hence, under the assumptions of the Lemma related to Assumptions B.1 and B.2, the distribution-pointwise inequalities from Theorem 3.1 in Mies and Steland [109] hold for each $P \in \mathcal{P}_n$. Since the suprema over all distributions in the collection \mathcal{P}_n of the upper bounds are finite, the distribution-uniform inequalities from the Lemma hold for \mathcal{P}_n by basic properties of the supremum. \square

Recently, Bonnerjee, Karmakar, and Wu [20] introduced univariate strong Gaussian approximation results with optimal rates and explicit constructions, building on prior work by Karmakar and Wu [81]. We emphasize that the distribution-uniform strong Gaussian approximation for high-dimensional nonstationary nonlinear processes from Lemma B.1 does *not* achieve this optimal rate. However, the convergence rates for the prediction errors from the estimation of the time-varying regression functions dominate the strong Gaussian approximation rates. Therefore, we do not “lose anything” by using Lemma B.1 in our regression-based conditional independence test instead of a distribution-uniform version of the strong Gaussian approximation from Bonnerjee, Karmakar, and Wu [20], as our main results would not change in any meaningful way.

The following result is a distribution-uniform version of Theorem 3.2 from Mies and Steland [109].

Lemma B.2. For some sample size $n \in \mathbb{N}$ and collection of distributions \mathcal{P}_n for the stochastic system $(G_{t,n}(\mathcal{H}_s))_{t \in [n], s \in \mathbb{Z}}$, let the \mathbb{R}^{d_n} -valued process $(W_{t,n})_{t \in [n]}$ be defined on the collection of probability spaces $(\Omega, \mathcal{B}, \mathbb{P}_P)_{P \in \mathcal{P}_n}$ so that $W_{t,n} = G_{t,n}(\mathcal{H}_t)$ with $W_{t,n} \in L_q(P)$ and $\theta_{P,t,n}(j, q, r)$ as in Definition B.1 for each $P \in \mathcal{P}_n$ and some $2 \leq r \leq q < \infty$. There exists a universal constant $K = K(q, r)$ such that

$$\begin{aligned} \sup_{P \in \mathcal{P}_n} \left(\mathbb{E}_P \max_{k \leq n} \left\| \sum_{t=1}^k (W_{t,n} - \mathbb{E}_P(W_{t,n})) \right\|_r^q \right)^{\frac{1}{q}} &\leq \sup_{P \in \mathcal{P}_n} \left(K n^{\frac{1}{2} - \frac{1}{q}} \sum_{j=1}^{\infty} \left(\sum_{t=1}^n \theta_{P,t,n}^q(j, q, r) \right)^{\frac{1}{q}} \right) \\ &\leq \sup_{P \in \mathcal{P}_n} \left(K n^{\frac{1}{2}} \sum_{j=1}^{\infty} \max_{t \leq n} \theta_{P,t,n}(j, q, r) \right). \end{aligned}$$

In the special case $r = 2$, the inequality may be improved to

$$\begin{aligned}
& \sup_{P \in \mathcal{P}_n} \left(\mathbb{E}_P \max_{k \leq n} \left\| \sum_{t=1}^k (W_{t,n} - \mathbb{E}_P(W_{t,n})) \right\|_2^q \right)^{\frac{1}{q}} \\
& \leq \sup_{P \in \mathcal{P}_n} \left(K \sum_{j=1}^{\infty} (j \wedge n)^{\frac{1}{2} - \frac{1}{q}} \left(\sum_{t=1}^n \theta_{P,t,n}^q(j, q, 2) \right)^{\frac{1}{q}} + K \sum_{j=1}^n \left(\sum_{t=1}^n \theta_{P,t,n}^2(j, 2, 2) \right)^{\frac{1}{2}} \right).
\end{aligned}$$

Proof of Lemma B.2: Under the assumptions of the Lemma, the distribution-pointwise inequalities from Theorem 3.2 in Mies and Steland [109] hold for each $P \in \mathcal{P}_n$. Since the suprema over all distributions in the collection \mathcal{P}_n of the upper bounds are always finite, the distribution-uniform inequalities from the Lemma hold for \mathcal{P}_n by basic properties of the supremum. \square

B.3 Distribution-uniform feasible Gaussian approximation

In this subsection, we introduce distribution-uniform versions of Theorem 4.1 and Proposition 4.2 from Mies and Steland [109] so that the distribution-uniform strong Gaussian approximation from Subsection B.2 can be used for statistical inference. The key is a distribution-uniform cumulative covariance estimator $\hat{Q}_{k,n}$ of the cumulative covariance matrices $Q_{P,k,n} = \sum_{t=1}^k \Sigma_{P,t,n}$ where $\Sigma_{P,t,n} = \sum_{h=-\infty}^{\infty} \text{Cov}_P(G_{t,n}(\mathcal{H}_0), G_{t,n}(\mathcal{H}_h))$ and $W_{t,n} = G_{t,n}(\mathcal{H}_t)$. We will prove these guarantees for the same estimator from Mies and Steland [109], namely

$$\hat{Q}_{k,n} = \sum_{r=L_n}^k \frac{1}{L_n} \left(\sum_{s=r-L_n+1}^r W_{s,n} \right)^{\otimes 2}$$

for some window size $L_n \asymp n^\zeta$ for some $\zeta \in (0, \frac{1}{2})$.

The following result is a distribution-uniform version of Theorem 4.1 from Mies and Steland [109].

Lemma B.3. *For some sample size $n \in \mathbb{N}$ and collection of distributions \mathcal{P}_n for the stochastic system $(G_{t,n}(\mathcal{H}_s))_{t \in [n], s \in \mathbb{Z}}$, let the \mathbb{R}^{d_n} -valued process $(W_{t,n})_{t \in [n]}$ be defined on the collection of probability spaces $(\Omega, \mathcal{B}, \mathbb{P}_P)_{P \in \mathcal{P}_n}$ so that $W_{t,n} = G_{t,n}(\mathcal{H}_t)$ and Assumptions B.1 and B.2 are satisfied for \mathcal{P}_n with $q \geq 4$ and $\beta > 2$. Then*

$$\sup_{P \in \mathcal{P}_n} \left(\mathbb{E}_P \max_{k=L_n, \dots, n} \left\| \hat{Q}_{k,n} - \sum_{t=1}^k \Sigma_{P,t,n} \right\|_{\text{tr}} \right) \leq K \Theta_n^2 \left(\Gamma_n \sqrt{L_n} + \sqrt{nd_n L_n} + nL_n^{-1} + nL_n^{2-\beta} \right)$$

for some universal constant K depending only on β and q .

Proof of Lemma B.3: Assumptions B.1 and B.2 are distribution-uniform versions of conditions (G.1) and (G.2) from Mies and Steland [109]. Hence, under the assumptions of the Lemma related to Assumptions B.1 and B.2, the distribution-pointwise inequalities from Theorem 4.1 in Mies and Steland [109] hold for each $P \in \mathcal{P}_n$. Since the supremum over all distributions in the collection \mathcal{P}_n of the upper bound is always finite, the distribution-uniform inequality from the Lemma holds for \mathcal{P}_n by basic properties of the supremum. \square

The next result is a distribution-uniform version of Proposition 4.2 from Mies and Steland [109].

Lemma B.4. *For some sample size $n \in \mathbb{N}$, let \mathcal{P}_n be a collection of distributions for the stochastic system $(G_{t,n}(\mathcal{H}_s))_{t \in [n], s \in \mathbb{Z}}$. Let $\Sigma_{P,t,n}, \Sigma'_{P,t,n} \in \mathbb{R}^{d_n \times d_n}$ be symmetric, positive definite matrices for each $t \in [n]$, $P \in \mathcal{P}_n$, and let $(V_{t,n})_{t \in [n]}$ be independent random vectors defined on the collection of probability spaces $(\Omega, \mathcal{B}, \mathbb{P}_P)_{P \in \mathcal{P}_n}$ so that $V_{t,n} \sim \mathcal{N}(0, \Sigma_{P,t,n})$ for each $t \in [n]$, $P \in \mathcal{P}_n$. Then, on a potentially enriched collection of probability spaces $(\Omega', \mathcal{B}', \mathbb{P}'_P)_{P \in \mathcal{P}_n}$, there exist independent random vectors $(V'_{t,n})_{t \in [n]}$ with $V'_{t,n} \sim \mathcal{N}(0, \Sigma'_{P,t,n})$ for each $t \in [n]$, $P \in \mathcal{P}_n$ such that*

$$\sup_{P \in \mathcal{P}_n} \left(\mathbb{E}_P \max_{k \leq n} \left\| \sum_{t=1}^k V_{t,n} - \sum_{t=1}^k V'_{t,n} \right\|_2^2 \right) \leq \sup_{P \in \mathcal{P}_n} \left(K \log(n) [\sqrt{n \delta_{P,n} \rho_{P,n}} + \rho_{P,n}] \right),$$

where

$$\begin{aligned}\delta_{P,n} &= \max_{k \leq n} \left\| \sum_{t=1}^k \Sigma_{P,t,n} - \sum_{t=1}^k \Sigma'_{P,t,n} \right\|_{\text{tr}}, \\ \rho_{P,n} &= \max_{t \leq n} \|\Sigma_{P,t,n}\|_{\text{tr}}.\end{aligned}$$

Proof of Lemma B.4: The distribution-pointwise inequalities from Proposition 4.2 in Mies and Steland [109] hold for each $P \in \mathcal{P}_n$. Since the supremum over all distributions in the collection \mathcal{P}_n of the upper bound is always finite, the distribution-uniform inequality from the Lemma holds for \mathcal{P}_n by basic properties of the supremum. \square

B.4 Auxiliary Lemmas

The following result is a distribution-uniform version of Proposition 5.4 from Mies and Steland [109].

Lemma B.5. *For some sample size $n \in \mathbb{N}$ and collection of distributions \mathcal{P}_n for the stochastic system $(G_{t,n}(\mathcal{H}_s))_{t \in [n], s \in \mathbb{Z}}$, let Assumption B.1 be satisfied for \mathcal{P}_n with some $q \geq 2$, $\beta > 0$, and constant $\Theta_n > 0$. Denote*

$$\gamma_{P,t,n}(h) = \text{Cov}_P[G_{t,n}(\mathcal{H}_0), G_{t,n}(\mathcal{H}_h)] \in \mathbb{R}^{d_n \times d_n}.$$

Then for all $t \in [n]$, $h \in \mathbb{Z}$, we have

$$\sup_{P \in \mathcal{P}_n} \|\gamma_{P,t,n}(h)\|_{\text{tr}} \leq \Theta_n^2 \sum_{j=h}^{\infty} j^{-\beta},$$

where $\|\cdot\|_{\text{tr}}$ denotes the trace norm. Hence, if $\beta > 2$, then the long-run covariance matrix

$$\gamma_{P,t,n} = \sum_{h=-\infty}^{\infty} \gamma_{P,t,n}(h),$$

is well-defined for all $t \in [n]$, $P \in \mathcal{P}_n$.

Proof of Lemma B.5: The distribution-pointwise inequality from Proposition 5.4 in Mies and Steland [109] holds for each $P \in \mathcal{P}_n$. Since the supremum over all distributions in the collection \mathcal{P}_n of the upper bound is always finite, the distribution-uniform inequality from the Lemma holds for \mathcal{P}_n by basic properties of the supremum. \square

The following result is a distribution-uniform version of the Rosenthal inequality from the first part of Theorem 5.6 from Mies and Steland [109].

Lemma B.6. *For some sample size $n \in \mathbb{N}$ and collection of distributions \mathcal{P}_n , let $(M_{t,n})_{t \in [n]}$ be a \mathbb{R}^{d_n} -valued martingale-difference sequence with distribution determined by $P \in \mathcal{P}_n$. For each $2 \leq r \leq q < \infty$, there exists a finite factor $C_{q,r}$ such that for any $n, d_n \in \mathbb{N}$, we have*

$$\begin{aligned}\sup_{P \in \mathcal{P}_n} \left(\mathbb{E}_P \max_{k \leq n} \left\| \sum_{t=1}^k M_{t,n} \right\|_r^q \right)^{\frac{1}{q}} &\leq C_{q,r} n^{\frac{1}{2} - \frac{1}{q}} \sup_{P \in \mathcal{P}_n} \left(\sum_{t=1}^n \mathbb{E}_P \|M_{t,n}\|_r^q \right)^{\frac{1}{q}} \\ &\leq C_{q,r} n^{\frac{1}{2}} \sup_{P \in \mathcal{P}_n} \left(\max_{t \leq n} (\mathbb{E}_P \|M_{t,n}\|_r^q)^{\frac{1}{q}} \right).\end{aligned}$$

Proof of Lemma B.6: The distribution-pointwise inequalities from the first part of Theorem 5.6 in Mies and Steland [109] hold for each $P \in \mathcal{P}_n$. Since the suprema over all distributions in the collection \mathcal{P}_n of the upper bounds are always finite, the distribution-uniform inequality from the Lemma holds for \mathcal{P}_n by basic properties of the supremum. \square

The following result is similar to the bounded convergence lemma from Lemma 25 in Shah and Peters [149].

Lemma B.7. *For some sample size $n \in \mathbb{N}$ and collection of distributions \mathcal{P}_n , let X_n be a generic real-valued random variable with distribution determined by $P \in \mathcal{P}_n$, where the collection of distributions \mathcal{P}_n can change with n . Let $K > 0$, and suppose that $|X_n| \leq K$ for all $n \in \mathbb{N}$ and $X_n = o_{\mathcal{P}}(1)$. Then we have*

$$\sup_{P \in \mathcal{P}_n} \mathbb{E}_P(|X_n|) = o(1).$$

Proof of Lemma B.7: For any given $\epsilon > 0$,

$$|X_n| = |X_n| \mathbb{1}_{\{|X_n| > \epsilon\}} + |X_n| \mathbb{1}_{\{|X_n| \leq \epsilon\}} \leq K \mathbb{1}_{\{|X_n| > \epsilon\}} + \epsilon.$$

By the assumption that $X_n = o_{\mathcal{P}}(1)$, we can find some $N \in \mathbb{N}$ such that $\sup_{P \in \mathcal{P}_n} \mathbb{P}_P(|X_n| > \epsilon) < \epsilon/K$ for $n \geq N$. Hence, for $n \geq N$ we have

$$\sup_{P \in \mathcal{P}_n} \mathbb{E}_P(|X_n|) \leq K \sup_{P \in \mathcal{P}_n} \mathbb{P}_P(|X_n| > \epsilon) + \epsilon < 2\epsilon.$$

Since $\epsilon > 0$ was arbitrary, we obtain the desired result. \square

C Proofs of Theoretical Results for dGCM and Sieve-dGCM

We will denote the three bias terms by

$$\begin{aligned}\hat{\mathbf{w}}_{P,t,n}^{\mathbf{f},\mathbf{g}} &= (\hat{w}_{P,t,n,m}^{\mathbf{f},\mathbf{g}})_{m \in \mathcal{D}_n} = (\hat{w}_{P,t,n,i,a}^{\mathbf{f}} \hat{w}_{P,t,n,j,b}^{\mathbf{g}})_{m \in \mathcal{D}_n}, \\ \hat{\mathbf{w}}_{P,t,n}^{\mathbf{g},\boldsymbol{\varepsilon}} &= (\hat{w}_{P,t,n,m}^{\mathbf{g},\boldsymbol{\varepsilon}})_{m \in \mathcal{D}_n} = (\hat{w}_{P,t,n,j,b}^{\mathbf{g}} \varepsilon_{P,t,n,i,a})_{m \in \mathcal{D}_n}, \\ \hat{\mathbf{w}}_{P,t,n}^{\mathbf{f},\boldsymbol{\xi}} &= (\hat{w}_{P,t,n,m}^{\mathbf{f},\boldsymbol{\xi}})_{m \in \mathcal{D}_n} = (\hat{w}_{P,t,n,i,a}^{\mathbf{f}} \xi_{P,t,n,j,b})_{m \in \mathcal{D}_n},\end{aligned}$$

where $m = (i, j, a, b) \in \mathcal{D}_n$. Also, denote

$$\begin{aligned}\hat{\mathbf{w}}_{P,n}^{\mathbf{f},\mathbf{g}} &= (\hat{\mathbf{w}}_{P,t,n}^{\mathbf{f},\mathbf{g}})_{t \in \mathcal{T}_{n,L}}, \\ \hat{\mathbf{w}}_{P,n}^{\mathbf{g},\boldsymbol{\varepsilon}} &= (\hat{\mathbf{w}}_{P,t,n}^{\mathbf{g},\boldsymbol{\varepsilon}})_{t \in \mathcal{T}_{n,L}}, \\ \hat{\mathbf{w}}_{P,n}^{\mathbf{f},\boldsymbol{\xi}} &= (\hat{\mathbf{w}}_{P,t,n}^{\mathbf{f},\boldsymbol{\xi}})_{t \in \mathcal{T}_{n,L}}.\end{aligned}$$

Note that when we write $o_{\mathcal{P}}(\cdot)$ and $O_{\mathcal{P}}(\cdot)$, we will always be doing so with reference to the collection of distributions $\mathcal{P}_{0,n}^*$ defined in the statement of the theorem.

C.1 Proof of Theorem 3.1

Step 1 (Bias Terms): We decompose the products of residuals into the products of errors and the three bias terms, and then apply the triangle inequality and subadditivity, which yields

$$\begin{aligned}& \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(S_{n,p}(\hat{\mathbf{R}}_n) > \hat{q}_{1-\alpha+\nu_n} + \tau_n) \\ & \leq \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(S_{n,p}(\mathbf{R}_{P,n}) > \hat{q}_{1-\alpha+\nu_n} + \frac{\tau_n}{2}) \\ & \quad + \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(S_{n,p}(\hat{\mathbf{w}}_{P,n}^{\mathbf{f},\mathbf{g}}) > \frac{\tau_n}{6}) \\ & \quad + \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(S_{n,p}(\hat{\mathbf{w}}_{P,n}^{\mathbf{g},\boldsymbol{\varepsilon}}) > \frac{\tau_n}{6}) \\ & \quad + \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(S_{n,p}(\hat{\mathbf{w}}_{P,n}^{\mathbf{f},\boldsymbol{\xi}}) > \frac{\tau_n}{6}).\end{aligned}$$

We will handle each of the three bias terms separately.

Step 1.1: Observe that for any $\delta > 0$, we have

$$\begin{aligned}& \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(\tau_n^{-1} S_{n,p}(\hat{\mathbf{w}}_{P,n}^{\mathbf{f},\mathbf{g}}) > \delta) \\ & \stackrel{(1)}{\leq} \delta^{-1} \tau_n^{-1} T_{n,L}^{-\frac{1}{2}} \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{E}_P \left(\max_{s \in \mathcal{T}_{n,L}} \left\| \sum_{t \leq s} \hat{\mathbf{w}}_{P,t,n}^{\mathbf{f},\mathbf{g}} \right\|_2 \right) \\ & \stackrel{(2)}{\leq} \delta^{-1} \tau_n^{-1} T_{n,L}^{-\frac{1}{2}} D_n \sup_{P \in \mathcal{P}_{0,n}^*} \max_{(i,j,a,b) \in \mathcal{D}_n} \mathbb{E}_P \left(\sum_{t \in \mathcal{T}_{n,L}} |\hat{w}_{P,t,n,i,a}^{\mathbf{f}}| |\hat{w}_{P,t,n,j,b}^{\mathbf{g}}| \right) \\ & \stackrel{(3)}{\leq} \delta^{-1} \tau_n^{-1} T_{n,L}^{\frac{1}{2}} D_n \sup_{P \in \mathcal{P}_{0,n}^*} \max_{(i,j,a,b) \in \mathcal{D}_n} \mathbb{E}_P \left(|\hat{w}_{P,t,n,i,a}^{\mathbf{f}}|^2 \right)^{\frac{1}{2}} \mathbb{E}_P \left(|\hat{w}_{P,t,n,j,b}^{\mathbf{g}}|^2 \right)^{\frac{1}{2}} \\ & \stackrel{(4)}{=} o(1),\end{aligned}$$

where the previous lines follow by (1) Markov's inequality and ℓ_p -norm inequalities, (2) the triangle inequality, ℓ_p -norm inequalities, linearity of expectation, (3) linearity of expectation and the Cauchy-Schwarz inequality, (4) the convergence rate requirements for the time-varying regression estimators.

Step 1.2: Observe that for any $\delta > 0$, we have

$$\begin{aligned}
& \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(\tau_n^{-1} S_{n,p}(\hat{\mathbf{w}}_{P,n}^{g,\varepsilon}) > \delta) \\
& \stackrel{(1)}{\leq} \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P \left(\tau_n^{-2} \max_{s \in \mathcal{T}_{n,L}} \left\| \frac{1}{\sqrt{T_{n,L}}} \sum_{t \leq s} \hat{\mathbf{w}}_{P,t,n}^{g,\varepsilon} \right\|_2^2 \geq \delta^2 \right) \\
& \stackrel{(2)}{\leq} \delta^{-2} \tau_n^{-2} T_{n,L}^{-1} \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{E}_P \left(\max_{s \in \mathcal{T}_{n,L}} \left\| \sum_{t \leq s} \hat{\mathbf{w}}_{P,t,n}^{g,\varepsilon} \right\|_2^2 \right) \\
& \stackrel{(3)}{\leq} \delta^{-2} \tau_n^{-2} T_{n,L}^{-1} (\bar{K} T_{n,L}^{\frac{1}{2}} D_n^{\frac{1}{2}} \sup_{P \in \mathcal{P}_{0,n}^*} \max_{t \in \mathcal{T}_{n,L}} \max_{j \in [d_Y]} \max_{b \in B_j} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g|^2)^{\frac{1}{2}})^2 \\
& \stackrel{(4)}{\leq} \delta^{-2} \tau_n^{-2} \bar{K}^2 D_n \sup_{P \in \mathcal{P}_{0,n}^*} \max_{t \in \mathcal{T}_{n,L}} \max_{j \in [d_Y]} \max_{b \in B_j} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g|^2) \\
& \stackrel{(5)}{=} o(1),
\end{aligned}$$

where the previous lines follow by (1) the assumption about the form of the test statistic and squaring, (2) Markov's inequality and linearity of expectation, (3) for some constant $\bar{K} > 0$ by the arguments below, (4) simplifying the expression, and (5) the convergence rate requirements for the time-varying regression estimator.

The following arguments are to show (3). These arguments are based on the constructions used in the proof of Theorem 3.2 in Mies and Steland [109], which build on the proof techniques from Theorem 1 in Liu, Xiao, and Wu [97]. For each $t \in \mathcal{T}_{n,L}$ and $h \in \mathbb{N}_0$, let

$$\mathcal{F}_{t,h}^{\hat{\mathbf{w}}^{g,\varepsilon}} = \sigma(\eta_t^\varepsilon, \eta_{t-1}^\varepsilon, \dots, \eta_{t-h}^\varepsilon, \mathcal{H}_t^{\hat{g}}),$$

where the input η_t^ε is from (11) and the input sequence $\mathcal{H}_t^{\hat{g}}$ is defined following Assumption 3.3. For each $n \in \mathbb{N}$, $P \in \mathcal{P}_{0,n}^*$, $t \in \mathcal{T}_{n,L}$, and $h \in \mathbb{N}_0$ let

$$\begin{aligned}
\hat{S}_{P,t,n,h}^{g,\varepsilon} &= \sum_{k \leq t} \hat{\mathbf{w}}_{P,k,n,h}^{g,\varepsilon}, \\
\hat{\mathbf{w}}_{P,t,n,h}^{g,\varepsilon} &= \mathbb{E}_P(\hat{\mathbf{w}}_{P,t,n}^{g,\varepsilon} | \mathcal{F}_{t,h}^{\hat{\mathbf{w}}^{g,\varepsilon}}), \\
\hat{\mathbf{w}}_{P,t,n,-1}^{g,\varepsilon} &= \mathbb{E}_P(\hat{\mathbf{w}}_{P,t,n}^{g,\varepsilon} | \mathcal{H}_t^{\hat{g}}) = \mathbf{0},
\end{aligned}$$

almost surely, because for each $n \in \mathbb{N}$, $P \in \mathcal{P}_{0,n}^*$, $(i, j, a, b) \in \mathcal{D}_n$, $t \in \mathcal{T}_{n,L}$ we have

$$\mathbb{E}_P(\hat{w}_{P,t,n,j,b}^g \varepsilon_{P,t,n,i,a} | \mathcal{H}_t^{\hat{g}}) = \hat{w}_{P,t,n,j,b}^g \mathbb{E}_P(\varepsilon_{P,t,n,i,a} | \mathcal{H}_t^{\hat{g}}) = 0, \quad (19)$$

almost surely, by Assumptions 3.3 and 3.4. For each $n \in \mathbb{N}$, $P \in \mathcal{P}_{0,n}^*$, $t \in \mathcal{T}_{n,L}$, and $h \in \mathbb{N}_0$ we have

$$\mathbb{E}_P(\|\hat{\mathbf{w}}_{P,t,n,h}^{g,\varepsilon}\|_2^2) < \infty, \quad (20)$$

by linearity of expectation, the contraction property of conditional expectation, and Assumption 3.3. For each $n \in \mathbb{N}$, $P \in \mathcal{P}_{0,n}^*$, $t \in \mathcal{T}_{n,L}$, and $h \in \mathbb{N}_0$, by the tower property we have

$$\mathbb{E}_P(\hat{\mathbf{w}}_{P,t,n,h+1}^{g,\varepsilon} | \mathcal{F}_{t,h}^{\hat{\mathbf{w}}^{g,\varepsilon}}) = \hat{\mathbf{w}}_{P,t,n,h}^{g,\varepsilon},$$

almost surely. Hence, for each $n \in \mathbb{N}$, $P \in \mathcal{P}_{0,n}^*$, and $t \in \mathcal{T}_{n,L}$, $(\hat{\mathbf{w}}_{P,t,n,h}^{g,\varepsilon})_{h=0}^\infty$ is a martingale with respect to the filtration $(\mathcal{F}_{t,h}^{\hat{\mathbf{w}}^{g,\varepsilon}})_{h=0}^\infty$. The martingale convergence theorem (see e.g. Theorem 1.5 of [123]) ensures that for each $n \in \mathbb{N}$, $P \in \mathcal{P}_{0,n}^*$, $t \in \mathcal{T}_{n,L}$ there exists some random vector $\tilde{\mathbf{w}}_{P,t,n}^{g,\varepsilon}$ such that $\mathbb{E}_P\|\tilde{\mathbf{w}}_{P,t,n}^{g,\varepsilon} - \hat{\mathbf{w}}_{P,t,n,h}^{g,\varepsilon}\|_2^2 \rightarrow 0$ as $h \rightarrow \infty$. The measurability of $\mathbf{G}_{P,t,n}^{\hat{\mathbf{w}}^{g,\varepsilon}}$ with respect to the projection σ -algebra, in view of Assumptions 3.1, 3.2, 3.3, 3.4, ensures that $\tilde{\mathbf{w}}_{P,t,n}^{g,\varepsilon} = \hat{\mathbf{w}}_{P,t,n}^{g,\varepsilon}$. Thus, for each $t \in \mathcal{T}_{n,L}$ we have

$$\hat{S}_{P,t,n}^{g,\varepsilon} = \sum_{k \leq t} \hat{\mathbf{w}}_{P,k,n}^{g,\varepsilon} = \sum_{h=0}^\infty (\hat{S}_{P,t,n,h}^{g,\varepsilon} - \hat{S}_{P,t,n,h-1}^{g,\varepsilon}), \quad (21)$$

by telescoping. For each $n \in \mathbb{N}$, $P \in \mathcal{P}_{0,n}^*$, and $h \in \mathbb{N}_0$,

$$(\hat{\mathbf{w}}_{P, \mathbb{T}_n^+ - k, n, h}^{g, \varepsilon} - \hat{\mathbf{w}}_{P, \mathbb{T}_n^+ - k, n, h-1}^{g, \varepsilon})_{k=0}^{\mathbb{T}_n^+ - \mathbb{T}_n^- - L_n},$$

are martingale differences with respect to the filtration $(\mathcal{G}_{\mathbb{T}_n^+, k, h}^{\hat{\mathbf{w}}^{g, \varepsilon}})_{k=0}^{\mathbb{T}_n^+ - \mathbb{T}_n^- - L_n}$, where

$$\mathcal{G}_{\mathbb{T}_n^+, k, h}^{\hat{\mathbf{w}}^{g, \varepsilon}} = \sigma(\mathcal{H}_{\mathbb{T}_n^+ - k}^{\hat{\mathbf{w}}^g}, \eta_{\mathbb{T}_n^+ - k - h}^\varepsilon, \eta_{\mathbb{T}_n^+ - k - h + 1}^\varepsilon \dots),$$

because for any $n \in \mathbb{N}$, $P \in \mathcal{P}_{0,n}^*$, $h \in \mathbb{N}_0$ and $k = 0, 1, \dots$, we have

$$\begin{aligned} & \mathbb{E}_P(\hat{\mathbf{w}}_{P, \mathbb{T}_n^+ - k, n, h}^{g, \varepsilon} - \hat{\mathbf{w}}_{P, \mathbb{T}_n^+ - k, n, h-1}^{g, \varepsilon} | \mathcal{G}_{\mathbb{T}_n^+, k-1, h}^{\hat{\mathbf{w}}^{g, \varepsilon}}) \\ &= \mathbb{E}_P(\mathbb{E}_P(\hat{\mathbf{w}}_{P, \mathbb{T}_n^+ - k, n}^{g, \varepsilon} | \mathcal{F}_{\mathbb{T}_n^+ - k, h}^{\hat{\mathbf{w}}^{g, \varepsilon}}) | \mathcal{G}_{\mathbb{T}_n^+, k-1, h}^{\hat{\mathbf{w}}^{g, \varepsilon}}) \\ &= \mathbb{E}_P(\mathbb{E}_P(\hat{\mathbf{w}}_{P, \mathbb{T}_n^+ - k, n}^{g, \varepsilon} | \mathcal{F}_{\mathbb{T}_n^+ - k, h-1}^{\hat{\mathbf{w}}^{g, \varepsilon}}) | \mathcal{G}_{\mathbb{T}_n^+, k-1, h}^{\hat{\mathbf{w}}^{g, \varepsilon}}) \\ &= \mathbf{0}, \end{aligned}$$

almost surely, because for each $n \in \mathbb{N}$, $P \in \mathcal{P}_{0,n}^*$, $(i, j, a, b) \in \mathcal{D}_n$, $h \in \mathbb{N}_0$ and $k = 0, 1, \dots$, we have

$$\begin{aligned} & \mathbb{E}_P(\mathbb{E}_P(\hat{w}_{P, \mathbb{T}_n^+ - k, n, j, b}^g \varepsilon_{P, \mathbb{T}_n^+ - k, n, i, a} | \mathcal{H}_{\mathbb{T}_n^+ - k}^{\hat{\mathbf{w}}^g}, \eta_{\mathbb{T}_n^+ - k - h}^\varepsilon, \eta_{\mathbb{T}_n^+ - k - h + 1}^\varepsilon, \dots, \eta_{\mathbb{T}_n^+ - k}^\varepsilon) \\ & | \mathcal{H}_{\mathbb{T}_n^+ - k + 1}^{\hat{\mathbf{w}}^g}, \eta_{\mathbb{T}_n^+ - k - h + 1}^\varepsilon, \eta_{\mathbb{T}_n^+ - k - h + 2}^\varepsilon, \dots) \\ & - \mathbb{E}_P(\mathbb{E}_P(\hat{w}_{P, \mathbb{T}_n^+ - k, n, j, b}^g \varepsilon_{P, \mathbb{T}_n^+ - k, n, i, a} | \mathcal{H}_{\mathbb{T}_n^+ - k}^{\hat{\mathbf{w}}^g}, \eta_{\mathbb{T}_n^+ - k - h + 1}^\varepsilon, \eta_{\mathbb{T}_n^+ - k - h + 2}^\varepsilon, \dots, \eta_{\mathbb{T}_n^+ - k}^\varepsilon) \\ & | \mathcal{H}_{\mathbb{T}_n^+ - k + 1}^{\hat{\mathbf{w}}^g}, \eta_{\mathbb{T}_n^+ - k - h + 1}^\varepsilon, \eta_{\mathbb{T}_n^+ - k - h + 2}^\varepsilon, \dots) \\ & \stackrel{(1)}{=} \hat{w}_{P, \mathbb{T}_n^+ - k, n, j, b}^g \mathbb{E}_P(\mathbb{E}_P(\varepsilon_{P, \mathbb{T}_n^+ - k, n, i, a} | \mathcal{H}_{\mathbb{T}_n^+ - k}^{\hat{\mathbf{w}}^g}, \eta_{\mathbb{T}_n^+ - k - h}^\varepsilon, \eta_{\mathbb{T}_n^+ - k - h + 1}^\varepsilon, \dots, \eta_{\mathbb{T}_n^+ - k}^\varepsilon) \\ & | \mathcal{H}_{\mathbb{T}_n^+ - k + 1}^{\hat{\mathbf{w}}^g}, \eta_{\mathbb{T}_n^+ - k - h + 1}^\varepsilon, \eta_{\mathbb{T}_n^+ - k - h + 2}^\varepsilon, \dots, \eta_{\mathbb{T}_n^+ - k}^\varepsilon) \\ & - \hat{w}_{P, \mathbb{T}_n^+ - k, n, j, b}^g \mathbb{E}_P(\mathbb{E}_P(\varepsilon_{P, \mathbb{T}_n^+ - k, n, i, a} | \mathcal{H}_{\mathbb{T}_n^+ - k}^{\hat{\mathbf{w}}^g}, \eta_{\mathbb{T}_n^+ - k - h + 1}^\varepsilon, \eta_{\mathbb{T}_n^+ - k - h + 2}^\varepsilon, \dots, \eta_{\mathbb{T}_n^+ - k}^\varepsilon) \\ & | \mathcal{H}_{\mathbb{T}_n^+ - k + 1}^{\hat{\mathbf{w}}^g}, \eta_{\mathbb{T}_n^+ - k - h + 1}^\varepsilon, \eta_{\mathbb{T}_n^+ - k - h + 2}^\varepsilon, \dots, \eta_{\mathbb{T}_n^+ - k}^\varepsilon) \\ & \stackrel{(2)}{=} \hat{w}_{P, \mathbb{T}_n^+ - k, n, j, b}^g \mathbb{E}_P(\varepsilon_{P, \mathbb{T}_n^+ - k, n, i, a} | \mathcal{H}_{\mathbb{T}_n^+ - k}^{\hat{\mathbf{w}}^g}, \eta_{\mathbb{T}_n^+ - k - h + 1}^\varepsilon, \eta_{\mathbb{T}_n^+ - k - h + 2}^\varepsilon, \dots, \eta_{\mathbb{T}_n^+ - k}^\varepsilon) \\ & - \hat{w}_{P, \mathbb{T}_n^+ - k, n, j, b}^g \mathbb{E}_P(\varepsilon_{P, \mathbb{T}_n^+ - k, n, i, a} | \mathcal{H}_{\mathbb{T}_n^+ - k}^{\hat{\mathbf{w}}^g}, \eta_{\mathbb{T}_n^+ - k - h + 1}^\varepsilon, \eta_{\mathbb{T}_n^+ - k - h + 2}^\varepsilon, \dots, \eta_{\mathbb{T}_n^+ - k}^\varepsilon) \\ & \stackrel{(3)}{=} 0, \end{aligned}$$

almost surely, by (1) Assumption 3.3, (2) the tower property and measurability, and (3) subtraction. Also, $\mathbb{E}_P(\|\hat{\mathbf{w}}_{P, \mathbb{T}_n^+ - k, n, h}^{g, \varepsilon} - \hat{\mathbf{w}}_{P, \mathbb{T}_n^+ - k, n, h-1}^{g, \varepsilon}\|_2^2) < \infty$ by the triangle inequality, squaring, linearity of expectation, the Cauchy-Schwarz inequality, and the same arguments as (20) (i.e. linearity of expectation, the contraction property of conditional expectation, and Assumption 3.3).

Next, observe that for each $n \in \mathbb{N}$ and $h \in \mathbb{N}_0$, we have

$$\begin{aligned} & \sup_{P \in \mathcal{P}_{0,n}^*} (\mathbb{E}_P \max_{t \in \mathcal{T}_{n,L}} \|\hat{\mathbf{S}}_{P, t, n, h}^{g, \varepsilon} - \hat{\mathbf{S}}_{P, t, n, h-1}^{g, \varepsilon}\|_2^2)^{\frac{1}{2}} \\ & \stackrel{(1)}{\leq} \sup_{P \in \mathcal{P}_{0,n}^*} (\mathbb{E}_P \|\hat{\mathbf{S}}_{P, \mathbb{T}_n^+, n, h}^{g, \varepsilon} - \hat{\mathbf{S}}_{P, \mathbb{T}_n^+, n, h-1}^{g, \varepsilon}\|_2^2)^{\frac{1}{2}} \\ & + \sup_{P \in \mathcal{P}_{0,n}^*} (\mathbb{E}_P \max_{t \in \mathcal{T}_{n,L}} \|(\hat{\mathbf{S}}_{P, \mathbb{T}_n^+, n, h}^{g, \varepsilon} - \hat{\mathbf{S}}_{P, \mathbb{T}_n^+, n, h-1}^{g, \varepsilon}) - (\hat{\mathbf{S}}_{P, t, n, h}^{g, \varepsilon} - \hat{\mathbf{S}}_{P, t, n, h-1}^{g, \varepsilon})\|_2^2)^{\frac{1}{2}} \\ & \stackrel{(2)}{\leq} \sup_{P \in \mathcal{P}_{0,n}^*} (\mathbb{E}_P \|\hat{\mathbf{S}}_{P, \mathbb{T}_n^+, n, h}^{g, \varepsilon} - \hat{\mathbf{S}}_{P, \mathbb{T}_n^+, n, h-1}^{g, \varepsilon}\|_2^2)^{\frac{1}{2}} \\ & + \sup_{P \in \mathcal{P}_{0,n}^*} \left(\mathbb{E}_P \max_{\ell=0, \dots, \mathbb{T}_n^+ - \mathbb{T}_n^- - L_n} \left\| \sum_{k=0}^{\ell} (\hat{\mathbf{w}}_{P, \mathbb{T}_n^+ - k, n, h}^{g, \varepsilon} - \hat{\mathbf{w}}_{P, \mathbb{T}_n^+ - k, n, h-1}^{g, \varepsilon}) \right\|_2^2 \right)^{\frac{1}{2}} \end{aligned}$$

by (1) adding and subtracting $\hat{\mathbf{S}}_{P, \mathbb{T}_n^+, n, h}^{g, \varepsilon} - \hat{\mathbf{S}}_{P, \mathbb{T}_n^+, n, h-1}^{g, \varepsilon}$ and the triangle inequality, (2) including the “last” term in this reversed partial sum and rewriting as the corresponding martingale. Continuing on from (2), we have

$$\begin{aligned}
& \sup_{P \in \mathcal{P}_{0,n}^*} (\mathbb{E}_P \|\hat{\mathbf{S}}_{P, \mathbb{T}_n^+, n, h}^{g, \varepsilon} - \hat{\mathbf{S}}_{P, \mathbb{T}_n^+, n, h-1}^{g, \varepsilon}\|_2^2)^{\frac{1}{2}} \\
& + \sup_{P \in \mathcal{P}_{0,n}^*} \left(\mathbb{E}_P \max_{\ell=0, \dots, \mathbb{T}_n^+ - \mathbb{T}_n^- - L_n} \left\| \sum_{k=0}^{\ell} (\hat{\mathbf{w}}_{P, \mathbb{T}_n^+ - k, n, h}^{g, \varepsilon} - \hat{\mathbf{w}}_{P, \mathbb{T}_n^+ - k, n, h-1}^{g, \varepsilon}) \right\|_2^2 \right)^{\frac{1}{2}} \\
& \stackrel{(3)}{\leq} 3 \sup_{P \in \mathcal{P}_{0,n}^*} (\mathbb{E}_P \|\hat{\mathbf{S}}_{P, \mathbb{T}_n^+, n, h}^{g, \varepsilon} - \hat{\mathbf{S}}_{P, \mathbb{T}_n^+, n, h-1}^{g, \varepsilon}\|_2^2)^{\frac{1}{2}} \\
& \stackrel{(4)}{\leq} K \sup_{P \in \mathcal{P}_{0,n}^*} \left(\sum_{t \in \mathcal{T}_{n,L}} \mathbb{E}_P \|\hat{\mathbf{w}}_{P, t, n, h}^{g, \varepsilon} - \hat{\mathbf{w}}_{P, t, n, h-1}^{g, \varepsilon}\|_2^2 \right)^{\frac{1}{2}},
\end{aligned}$$

by (3) Doob’s maximal inequality (see e.g. Theorem 1.9 of [123]), and (4) upper bounding by max of partial sums and applying Lemma B.6 with the finite constant $K/3 > 0$. Hence, we have the inequality

$$\begin{aligned}
& \sup_{P \in \mathcal{P}_{0,n}^*} (\mathbb{E}_P \max_{t \in \mathcal{T}_{n,L}} \|\hat{\mathbf{S}}_{P, t, n, h}^{g, \varepsilon} - \hat{\mathbf{S}}_{P, t, n, h-1}^{g, \varepsilon}\|_2^2)^{\frac{1}{2}} \\
& \leq K \sup_{P \in \mathcal{P}_{0,n}^*} \left(\sum_{t \in \mathcal{T}_{n,L}} \mathbb{E}_P \|\hat{\mathbf{w}}_{P, t, n, h}^{g, \varepsilon} - \hat{\mathbf{w}}_{P, t, n, h-1}^{g, \varepsilon}\|_2^2 \right)^{\frac{1}{2}}.
\end{aligned} \tag{22}$$

Observe that for $h = 1, 2, \dots$, we have

$$\begin{aligned}
& \mathbb{E}_P \|\hat{\mathbf{w}}_{P, t, n, h}^{g, \varepsilon} - \hat{\mathbf{w}}_{P, t, n, h-1}^{g, \varepsilon}\|_2^2 \\
& \stackrel{(1)}{=} \sum_{m=(i,j,a,b) \in \mathcal{D}_n} \mathbb{E}_P (|\mathbb{E}_P(\hat{w}_{P, t, n, j, b}^g \varepsilon_{P, t, n, i, a} | \eta_t^\varepsilon, \dots, \eta_{t-h}^\varepsilon, \mathcal{H}_t^{\hat{g}}) \\
& - \mathbb{E}_P(\hat{w}_{P, t, n, j, b}^g \varepsilon_{P, t, n, i, a} | \eta_t^\varepsilon, \dots, \eta_{t-h+1}^\varepsilon, \mathcal{H}_t^{\hat{g}})|^2) \\
& \stackrel{(2)}{=} \sum_{m=(i,j,a,b) \in \mathcal{D}_n} \mathbb{E}_P (|\mathbb{E}_P(\hat{w}_{P, t, n, j, b}^g (\varepsilon_{P, t, n, i, a} | \eta_t^\varepsilon, \dots, \eta_{t-h}^\varepsilon, \mathcal{H}_t^{\hat{g}}) \\
& - \mathbb{E}_P(\varepsilon_{P, t, n, i, a} | \eta_t^\varepsilon, \dots, \eta_{t-h+1}^\varepsilon, \mathcal{H}_t^{\hat{g}}))|^2) \\
& \stackrel{(3)}{=} \sum_{m=(i,j,a,b) \in \mathcal{D}_n} \mathbb{E}_P (|\hat{w}_{P, t, n, j, b}^g \mathbb{E}_P[(\mathbb{E}_P(\varepsilon_{P, t, n, i, a} | \eta_{t,a}^\varepsilon, \dots, \eta_{t-h,a}^\varepsilon, \mathcal{H}_t^{\hat{g}}) \\
& - \mathbb{E}_P(\varepsilon_{P, t, n, i, a} | \eta_{t,a}^\varepsilon, \dots, \eta_{t-h+1,a}^\varepsilon, \mathcal{H}_t^{\hat{g}})) | \eta_{t,a}^\varepsilon, \dots, \eta_{t-h,a}^\varepsilon, \mathcal{H}_t^{\hat{g}}]|^2) \\
& \stackrel{(4)}{=} \sum_{m=(i,j,a,b) \in \mathcal{D}_n} \mathbb{E}_P (|\hat{w}_{P, t, n, j, b}^g \mathbb{E}_P[(\mathbb{E}_P(G_{P, t, n, i, a}^\varepsilon(\mathcal{H}_{t,a}^\varepsilon) | \eta_{t,a}^\varepsilon, \dots, \eta_{t-h,a}^\varepsilon, \mathcal{H}_t^{\hat{g}}) \\
& - \mathbb{E}_P(G_{P, t, n, i, a}^\varepsilon(\mathcal{H}_{t,a}^\varepsilon) | \eta_{t,a}^\varepsilon, \dots, \eta_{t-h+1,a}^\varepsilon, \mathcal{H}_t^{\hat{g}})) | \eta_{t,a}^\varepsilon, \dots, \eta_{t-h,a}^\varepsilon, \mathcal{H}_t^{\hat{g}}]|^2) \\
& \stackrel{(5)}{=} \sum_{m=(i,j,a,b) \in \mathcal{D}_n} \mathbb{E}_P (|\hat{w}_{P, t, n, j, b}^g \mathbb{E}_P[(\mathbb{E}_P(G_{P, t, n, i, a}^\varepsilon(\mathcal{H}_{t,a}^\varepsilon) | \eta_{t,a}^\varepsilon, \dots, \eta_{t-h,a}^\varepsilon, \mathcal{H}_t^{\hat{g}}) \\
& - \mathbb{E}_P(G_{P, t, n, i, a}^\varepsilon(\tilde{\mathcal{H}}_{t,a,h}^\varepsilon) | \eta_{t,a}^\varepsilon, \dots, \eta_{t-h,a}^\varepsilon, \mathcal{H}_t^{\hat{g}})) | \eta_{t,a}^\varepsilon, \dots, \eta_{t-h,a}^\varepsilon, \mathcal{H}_t^{\hat{g}}]|^2),
\end{aligned}$$

by (1) rewriting the expression, (2) the causal representation from Assumption 3.3, (3) measurability of the conditional expectations and the linearity property of conditional expectation, (4) the causal representation from Assumption 3.4, and (5) replacing $\eta_{t-h,a}^\varepsilon$ with the iid copy $\tilde{\eta}_{t-h,a}^\varepsilon$. Continuing on

from line (5), we have

$$\begin{aligned}
& \sum_{m=(i,j,a,b) \in \mathcal{D}_n} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g| \mathbb{E}_P[(\mathbb{E}_P(G_{P,t,n,i,a}^\varepsilon(\mathcal{H}_{t,a}^\varepsilon)|\eta_{t,a}^\varepsilon, \dots, \eta_{t-h,a}^\varepsilon, \mathcal{H}_t^{\hat{g}}) \\
& - \mathbb{E}_P(G_{P,t,n,i,a}^\varepsilon(\tilde{\mathcal{H}}_{t,a,h}^\varepsilon)|\eta_{t,a}^\varepsilon, \dots, \eta_{t-h,a}^\varepsilon, \mathcal{H}_t^{\hat{g}}))|\eta_{t,a}^\varepsilon, \dots, \eta_{t-h,a}^\varepsilon, \mathcal{H}_t^{\hat{g}}]|^2) \\
& \stackrel{(6)}{=} \sum_{m=(i,j,a,b) \in \mathcal{D}_n} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g| \mathbb{E}_P[G_{P,t,n,i,a}^\varepsilon(\mathcal{H}_{t,a}^\varepsilon) - G_{P,t,n,i,a}^\varepsilon(\tilde{\mathcal{H}}_{t,a,h}^\varepsilon)|\eta_{t,a}^\varepsilon, \dots, \eta_{t-h,a}^\varepsilon, \mathcal{H}_t^{\hat{g}}]|^2) \\
& \stackrel{(7)}{\leq} D_n \max_{(i,j,a,b) \in \mathcal{D}_n} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g|^2) (\theta_{P,t,n,i,a}^{\varepsilon,\infty}(h))^2 \\
& \stackrel{(8)}{\leq} D_n \max_{(i,j,a,b) \in \mathcal{D}_n} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g|^2) (\bar{\Theta}^\infty(h \vee 1))^{-\bar{\beta}^\infty})^2,
\end{aligned}$$

by (6) measurability and linearity of the conditional expectations, and (7) Hölder's inequality, contraction property of conditional expectation, rewriting as the functional dependence measure from Definition 3.1, and upper bounding by the sum by D_n times the maximum over the dimension/time-offset combinations in \mathcal{D}_n , and (8) the upper bound on the L^∞ functional dependence measure from Assumption 3.5. Similarly, for $h = 0$, we have

$$\begin{aligned}
& \mathbb{E}_P \|\hat{w}_{P,t,n,0}^{g,\varepsilon} - \hat{w}_{P,t,n,-1}^{g,\varepsilon}\|_2^2 \\
& \stackrel{(1)}{=} \mathbb{E}_P \|\hat{w}_{P,t,n,0}^{g,\varepsilon}\|_2^2 \\
& \stackrel{(2)}{=} \sum_{m=(i,j,a,b) \in \mathcal{D}_n} \mathbb{E}_P(|\mathbb{E}_P(\hat{w}_{P,t,n,j,b}^g \varepsilon_{P,t,n,i,a} | \eta_t^\varepsilon, \mathcal{H}_t^{\hat{g}})|^2) \\
& \stackrel{(3)}{=} \sum_{m=(i,j,a,b) \in \mathcal{D}_n} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g| \mathbb{E}_P(\varepsilon_{P,t,n,i,a} | \eta_t^\varepsilon, \mathcal{H}_t^{\hat{g}})|^2) \\
& \stackrel{(4)}{\leq} D_n \max_{(i,j,a,b) \in \mathcal{D}_n} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g|^2) (\bar{\Theta}^\infty)^2,
\end{aligned}$$

because (1) $\hat{w}_{P,t,n,-1}^{g,\varepsilon} = 0$, (2) rewriting the expression, (3) Assumption 3.3, and (4) Hölder's inequality, contraction property of conditional expectation, applying the upper bound on the L^∞ norm from Assumption 3.5, and upper bounding by the sum by D_n times the maximum over the dimension/time-offset combinations in \mathcal{D}_n . Hence, for all $h \in \mathbb{N}_0$ we have

$$\mathbb{E}_P \|\hat{w}_{P,t,n,h}^{g,\varepsilon} - \hat{w}_{P,t,n,h-1}^{g,\varepsilon}\|_2^2 \leq D_n \max_{(i,j,a,b) \in \mathcal{D}_n} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g|^2) (\bar{\Theta}^\infty(h \vee 1))^{-\bar{\beta}^\infty})^2. \quad (23)$$

Summing over $h \in \mathbb{N}_0$, we have

$$\begin{aligned}
& \sup_{P \in \mathcal{P}_{0,n}^*} (\mathbb{E}_P \max_{t \in \mathcal{T}_{n,L}} \|\hat{S}_{P,t,n}^{g,\varepsilon}\|_2^2)^{\frac{1}{2}} \\
& \stackrel{(1)}{\leq} \sum_{h=0}^{\infty} \sup_{P \in \mathcal{P}_{0,n}^*} (\mathbb{E}_P \max_{t \in \mathcal{T}_{n,L}} \|\hat{S}_{P,t,n,h}^{g,\varepsilon} - \hat{S}_{P,t,n,h-1}^{g,\varepsilon}\|_2^2)^{\frac{1}{2}} \\
& \stackrel{(2)}{\leq} \sum_{h=0}^{\infty} K \sup_{P \in \mathcal{P}_{0,n}^*} \left(\sum_{t \in \mathcal{T}_{n,L}} \mathbb{E}_P \|\hat{w}_{P,t,n,h}^{g,\varepsilon} - \hat{w}_{P,t,n,h-1}^{g,\varepsilon}\|_2^2 \right)^{\frac{1}{2}} \\
& \stackrel{(3)}{\leq} \sum_{h=0}^{\infty} K \sup_{P \in \mathcal{P}_{0,n}^*} \left(\sum_{t \in \mathcal{T}_{n,L}} D_n \max_{j \in [d_Y]} \max_{b \in B_j} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g|^2) (\bar{\Theta}^\infty(h \vee 1))^{-\bar{\beta}^\infty})^2 \right)^{\frac{1}{2}},
\end{aligned}$$

by (1) the telescoping argument from (21) and the triangle inequality, (2) applying the inequality (22),

and (3) applying the inequality (23). Continuing on from line (3), we have

$$\begin{aligned}
& \sum_{h=0}^{\infty} K \sup_{P \in \mathcal{P}_{0,n}^*} \left(\sum_{t \in \mathcal{T}_{n,L}} D_n \max_{j \in [d_Y]} \max_{b \in B_j} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g|^2) (\bar{\Theta}^\infty (h \vee 1)^{-\bar{\beta}^\infty})^2 \right)^{\frac{1}{2}} \\
& \stackrel{(4)}{\leq} \sum_{h=0}^{\infty} K \sup_{P \in \mathcal{P}_{0,n}^*} (T_{n,L} D_n \max_{t \in \mathcal{T}_{n,L}} \max_{j \in [d_Y]} \max_{b \in B_j} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g|^2) (\bar{\Theta}^\infty (h \vee 1)^{-\bar{\beta}^\infty})^2)^{\frac{1}{2}} \\
& \stackrel{(5)}{\leq} \bar{\Theta}^\infty K T_{n,L}^{\frac{1}{2}} D_n^{\frac{1}{2}} \sup_{P \in \mathcal{P}_{0,n}^*} \max_{t \in \mathcal{T}_{n,L}} \max_{j \in [d_Y]} \max_{b \in B_j} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g|^2)^{\frac{1}{2}} \sum_{h=0}^{\infty} (h \vee 1)^{-\bar{\beta}^\infty} \\
& \stackrel{(6)}{\leq} \bar{\Theta}^\infty \bar{K}^\infty K T_{n,L}^{\frac{1}{2}} D_n^{\frac{1}{2}} \sup_{P \in \mathcal{P}_{0,n}^*} \max_{t \in \mathcal{T}_{n,L}} \max_{j \in [d_Y]} \max_{b \in B_j} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g|^2)^{\frac{1}{2}} \\
& \stackrel{(7)}{\leq} \bar{K} T_{n,L}^{\frac{1}{2}} D_n^{\frac{1}{2}} \sup_{P \in \mathcal{P}_{0,n}^*} \max_{t \in \mathcal{T}_{n,L}} \max_{j \in [d_Y]} \max_{b \in B_j} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g|^2)^{\frac{1}{2}},
\end{aligned}$$

by (4) upper bounding each term by the maximum over time t , (5) simplifying the expression, (6) writing $\bar{K}^\infty = \sum_{h=0}^{\infty} (h \vee 1)^{-\bar{\beta}^\infty} < \infty$ since $\bar{\beta}^\infty > 1$ by upper Assumption 3.5, and (7) grouping together the positive constants into the positive constant \bar{K} .

Step 1.3: The same arguments as Step 1.2 (i.e. exchanging g, ε with f, ξ) can be used to show that for $n \in \mathbb{N}$ and $\delta > 0$ we have

$$\sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(\tau_n^{-1} S_{n,p}(\hat{\mathbf{w}}_{P,n}^{f,\xi}) > \delta) = o(1).$$

Step 2 (Strong Gaussian Approximation): Next, we turn to the products of errors $(\mathbf{R}_{P,t,n})_{t \in \mathcal{T}_{n,L}}$. Denote the Gaussian random vectors associated with the strong Gaussian approximation of the product of errors by $\mathbf{R}_{t,n}^\dagger \sim \mathcal{N}(0, \Sigma_{P,t,n}^{\mathbf{R}})$ for $t \in \mathcal{T}_{n,L}$. Observe that

$$\begin{aligned}
& \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(S_{n,p}(\mathbf{R}_{P,n}) > \hat{q}_{1-\alpha+\nu_n} + \frac{\tau_n}{2}) \\
& \stackrel{(1)}{\leq} \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(S_{n,p}(\mathbf{R}_n^\dagger) > \hat{q}_{1-\alpha+\nu_n} + \frac{\tau_n}{4}) \\
& + \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P \left(\max_{s \in \mathcal{T}_{n,L}} \left\| \frac{1}{\sqrt{T_{n,L}}} \sum_{t \leq s} (\mathbf{R}_{P,t,n} - \mathbf{R}_{t,n}^\dagger) \right\|_2 > \frac{\tau_n}{4} \right) \\
& \stackrel{(2)}{\leq} \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(S_{n,p}(\mathbf{R}_n^\dagger) > \hat{q}_{1-\alpha+\nu_n} + \frac{\tau_n}{4}) \\
& + 4\tau_n^{-1} \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{E}_P \left(\max_{s \in \mathcal{T}_{n,L}} \left\| \frac{1}{\sqrt{T_{n,L}}} \sum_{t \leq s} (\mathbf{R}_{P,t,n} - \mathbf{R}_{t,n}^\dagger) \right\|_2 \right) \\
& \stackrel{(3)}{\leq} \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(S_{n,p}(\mathbf{R}_n^\dagger) > \hat{q}_{1-\alpha+\nu_n} + \frac{\tau_n}{4}) \\
& + 4\tau_n^{-1} K D_n^{\frac{1}{2}} \bar{\Theta}^R (\bar{\Gamma}_n^R)^{\frac{1}{2}} \frac{\bar{\beta}^R - 2}{\bar{\beta}^R - 1} \sqrt{\log(T_{n,L})} \left(\frac{D_n}{T_{n,L}} \right)^{\xi(\bar{q}^R, \bar{\beta}^R)},
\end{aligned}$$

where (1) follows from the triangle inequality, subadditivity, and the assumption about the form of the test statistic, (2) follows by Markov's inequality, and (3) follows by the distribution-uniform strong Gaussian approximation for high-dimensional nonstationary processes from Lemma B.1. By

subadditivity and monotonicity, we have

$$\begin{aligned}
& \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(S_{n,p}(\mathbf{R}_n^\dagger) > \hat{q}_{1-\alpha+\nu_n} + \frac{\tau_n}{4}) \\
& \leq \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(S_{n,p}(\mathbf{R}_n^\dagger) > q_{1-\alpha}) \\
& + \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(q_{1-\alpha} > \hat{q}_{1-\alpha+\nu_n} + \frac{\tau_n}{4}) \\
& = \alpha + \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(q_{1-\alpha} > \hat{q}_{1-\alpha+\nu_n} + \frac{\tau_n}{4}).
\end{aligned}$$

Step 3 (Covariance Approximation): Now, we focus on upper bounding

$$\sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(q_{1-\alpha} > \hat{q}_{1-\alpha+\nu_n} + \frac{\tau_n}{4}).$$

Step 3.1: Let us reflect on the implications of Proposition 4.2 of Mies and Steland [109], which is the distribution-pointwise version of Lemma B.4. Proposition 4.2 states that for each $n \in \mathbb{N}$ and $P \in \mathcal{P}_{0,n}^*$, for some cumulative covariance $\bar{Q}_{P,n}^{\mathbf{R}}$, there exist *independent* Gaussian random vectors $\bar{\mathbf{R}}_{t,n} \sim \mathcal{N}(0, \bar{\Sigma}_{P,t,n}^{\mathbf{R}})$ for $t \in \mathcal{T}_{n,L}$ with $\bar{\Sigma}_{P,t,n}^{\mathbf{R}} = \bar{Q}_{P,t,n}^{\mathbf{R}} - \bar{Q}_{P,t-1,n}^{\mathbf{R}}$ that are coupled with the Gaussian random vectors from the strong Gaussian approximation of the product of errors $\mathbf{R}_{t,n}^\dagger \sim \mathcal{N}(0, \Sigma_{P,t,n}^{\mathbf{R}})$ for $t \in \mathcal{T}_{n,L}$, such that

$$\mathbb{E}_P \max_{k \in \mathcal{T}_{n,L}} \left\| \sum_{t \leq k} \mathbf{R}_{t,n}^\dagger - \sum_{t \leq k} \bar{\mathbf{R}}_{t,n} \right\|_2^2 \leq K \log(T_{n,L}) [\sqrt{T_{n,L} \bar{\delta}_{P,n} \rho_{P,n}} + \rho_{P,n}] = \bar{\Delta}_{P,n},$$

where

$$\bar{\delta}_{P,n} = \max_{k \in \mathcal{T}_{n,L}} \left\| \sum_{t \leq k} \Sigma_{P,t,n}^{\mathbf{R}} - \sum_{t \leq k} \bar{\Sigma}_{P,t,n}^{\mathbf{R}} \right\|_{\text{tr}}$$

and

$$\rho_{P,n} = \max_{t \in \mathcal{T}_{n,L}} \|\Sigma_{P,t,n}^{\mathbf{R}}\|_{\text{tr}}.$$

Let $\bar{\mathbf{R}}_n = (\bar{\mathbf{R}}_{t,n})_{t \in \mathcal{T}_{n,L}}$ and denote the $(1 - \alpha)$ quantile of $S_{n,p}(\bar{\mathbf{R}}_n)$ by $\bar{q}_{1-\alpha}$. For each $n \in \mathbb{N}$ and $P \in \mathcal{P}_{0,n}^*$, we have

$$\begin{aligned}
& \mathbb{P}_P(S_{n,p}(\mathbf{R}_n^\dagger) > \bar{q}_{1-\alpha+\nu_n} + \frac{\tau_n}{4}) \\
& \stackrel{(1)}{\leq} \mathbb{P}_P(S_{n,p}(\bar{\mathbf{R}}_n) > \bar{q}_{1-\alpha+\nu_n}) \\
& + \mathbb{P}_P \left(\max_{s \in \mathcal{T}_{n,L}} \left\| \frac{1}{\sqrt{T_{n,L}}} \sum_{t \leq s} (\mathbf{R}_{t,n}^\dagger - \bar{\mathbf{R}}_{t,n}) \right\|_2 > \frac{\tau_n}{4} \right) \\
& \stackrel{(2)}{=} \mathbb{P}_P(S_{n,p}(\bar{\mathbf{R}}_n) > \bar{q}_{1-\alpha+\nu_n}) \\
& + \mathbb{P}_P \left(\max_{s \in \mathcal{T}_{n,L}} \left\| \frac{1}{\sqrt{T_{n,L}}} \sum_{t \leq s} (\mathbf{R}_{t,n}^\dagger - \bar{\mathbf{R}}_{t,n}) \right\|_2^2 > \frac{\tau_n^2}{16} \right) \\
& \stackrel{(3)}{\leq} \mathbb{P}_P(S_{n,p}(\bar{\mathbf{R}}_n) > \bar{q}_{1-\alpha+\nu_n}) \\
& + 16\tau_n^{-2} T_{n,L}^{-1} \mathbb{E}_P \left(\max_{s \in \mathcal{T}_{n,L}} \left\| \sum_{t \leq s} (\mathbf{R}_{t,n}^\dagger - \bar{\mathbf{R}}_{t,n}) \right\|_2^2 \right) \\
& \stackrel{(4)}{\leq} (\alpha - \nu_n) + 16\tau_n^{-2} \bar{\Delta}_{P,n} T_{n,L}^{-1} \stackrel{(5)}{=} \alpha + \left[16\tau_n^{-2} \bar{\Delta}_{P,n} T_{n,L}^{-1} - \nu_n \right],
\end{aligned}$$

where the previous lines follow by (1) the triangle inequality, subadditivity, the assumption about the form of the test statistic, (2) squaring, (3) Markov's inequality, (4) Proposition 4.2 from Mies and Steland [109], and (5) rearranging terms. We see that if

$$\left[16\tau_n^{-2}\bar{\Delta}_{P,n}T_{n,L}^{-1} - \nu_n\right] < 0,$$

then

$$\mathbb{P}_P(S_{n,p}(\mathbf{R}_n^\dagger) > \bar{q}_{1-\alpha+\nu_n} + \frac{\tau_n}{4}) < \alpha,$$

which implies that $\bar{q}_{1-\alpha+\nu_n} + \frac{\tau_n}{4}$ is greater than $q_{1-\alpha}^\dagger$, the $(1-\alpha)$ quantile of $S_{n,p}(\mathbf{R}_n^\dagger)$. Hence, if

$$q_{1-\alpha}^\dagger \geq \bar{q}_{1-\alpha+\nu_n} + \frac{\tau_n}{4},$$

then

$$\left[16\tau_n^{-2}\bar{\Delta}_{P,n}T_{n,L}^{-1} - \nu_n\right] \geq 0,$$

or equivalently

$$\bar{\Delta}_{P,n} \geq \frac{1}{16}T_{n,L}\nu_n\tau_n^2.$$

Step 3.2: Now, we apply this idea with the cumulative covariance of the residual products \hat{Q}_n^R . By the implication stated at the end of Step 3.1 and monotonicity, we have

$$\begin{aligned} & \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(q_{1-\alpha} > \hat{q}_{1-\alpha+\nu_n} + \frac{\tau_n}{4}) \\ & \leq \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(\hat{\Delta}_{P,n} \geq \frac{1}{16}T_{n,L}\nu_n\tau_n^2), \end{aligned}$$

where we have replaced $\bar{\Delta}_{P,n}$, $\bar{\delta}_{P,n}$ with $\hat{\Delta}_{P,n}$, $\hat{\delta}_{P,n}$ which are defined by

$$\hat{\Delta}_{P,n} = K \log(T_{n,L}) \left[\sqrt{T_{n,L}\hat{\delta}_{P,n}\rho_{P,n} + \rho_{P,n}} \right],$$

$$\hat{\delta}_{P,n} = \max_{k \in \mathcal{T}_{n,L}} \left\| \sum_{t \leq k} \boldsymbol{\Sigma}_{P,t,n}^R - \hat{Q}_{k,n}^R \right\|_{\text{tr}},$$

and $\rho_{P,n}$ is defined in the same way as

$$\rho_{P,n} = \max_{t \in \mathcal{T}_{n,L}} \|\boldsymbol{\Sigma}_{P,t,n}^R\|_{\text{tr}}.$$

Thus, if we can find φ_n such that $\hat{\Delta}_{P,n} = O_{\mathcal{P}}(\varphi_n)$ and if we select the offsets so that $\nu_n\tau_n^2 \gg T_{n,L}^{-1}\varphi_n$, or equivalently $\nu_n \gg \tau_n^{-2}T_{n,L}^{-1}\varphi_n$, then we will have

$$\sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(\hat{\Delta}_{P,n} \geq \frac{1}{16}T_{n,L}\nu_n\tau_n^2) \rightarrow 0.$$

By Lemma B.5 and Assumption 3.5, we have

$$\sup_{P \in \mathcal{P}_{0,n}^*} \rho_{P,n} \leq K_\rho D_n (\bar{\Theta}^R)^2,$$

for some constant $K_\rho > 0$, so we obtain $\hat{\Delta}_{P,n} = O_{\mathcal{P}}(\varphi_n)$ with

$$\varphi_n = \log(T_{n,L})D_n \left[T_{n,L}^{\frac{1}{2}}D_n^{-\frac{1}{2}}\hat{\delta}_{P,n}^{\frac{1}{2}} + 1 \right].$$

Plugging φ_n into the offset condition $\nu_n \gg \tau_n^{-2}T_{n,L}^{-1}\varphi_n$ that we wish to satisfy, if we have

$$\nu_n \gg \log(T_{n,L})D_n(\tau_n^{-2}(T_{n,L}^{-\frac{1}{2}}D_n^{-\frac{1}{2}}\hat{\delta}_{P,n}^{\frac{1}{2}} + T_{n,L}^{-1})),$$

then

$$\sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(\hat{\Delta}_{P,n} \geq \frac{1}{16} T_{n,L} \nu_n \tau_n^2) \rightarrow 0.$$

Step 3.3: It remains to analyze $\hat{\delta}_{P,n}$. By the triangle inequality, we have

$$\begin{aligned} \hat{\delta}_{P,n} &= \max_{k \in \mathcal{T}_{n,L}} \left\| \sum_{t \leq k} \Sigma_{P,t,n}^R - \hat{Q}_{k,n}^R \right\|_{\text{tr}} \\ &\leq \max_{k \in \mathcal{T}_{n,L}} \left\| \sum_{t \leq k} \Sigma_{P,t,n}^R - Q_{P,k,n}^R \right\|_{\text{tr}} \\ &\quad + \max_{k \in \mathcal{T}_{n,L}} \|\hat{Q}_{k,n}^R - Q_{P,k,n}^R\|_{\text{tr}}. \end{aligned}$$

By Lemma B.3, Assumption 3.5, and Assumption 3.6, the covariance estimation error can be bounded as

$$\begin{aligned} &\sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{E}_P \left(\max_{k \in \mathcal{T}_{n,L}} \left\| \sum_{t \leq k} \Sigma_{P,t,n}^R - Q_{P,k,n}^R \right\|_{\text{tr}} \right) \\ &\leq K(\bar{\Theta}^R)^2 D_n (\bar{\Gamma}_n^R L_n^{\frac{1}{2}} + T_{n,L}^{\frac{1}{2}} D_n^{\frac{1}{2}} L_n^{\frac{1}{2}} + T_{n,L} L_n^{-1} + T_{n,L} L_n^{2-\bar{\beta}^R}) \\ &= O(r_{n,1}^\delta), \end{aligned}$$

where

$$r_{n,1}^\delta = D_n (\bar{\Gamma}_n^R L_n^{\frac{1}{2}} + T_{n,L}^{\frac{1}{2}} D_n^{\frac{1}{2}} L_n^{\frac{1}{2}} + T_{n,L} L_n^{-1} + T_{n,L} L_n^{2-\bar{\beta}^R}).$$

Next, we must handle the prediction errors due using the residual products instead of the error products. For any $\epsilon > 0$, we have

$$\begin{aligned} &\sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{E}_P \left(\max_{k \in \mathcal{T}_{n,L}} \|\hat{Q}_{k,n}^R - Q_{P,k,n}^R\|_{\text{tr}} \wedge \epsilon \right) \\ &\stackrel{(1)}{=} \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{E}_P \left(\max_{k \in \mathcal{T}_{n,L}} \left\| \frac{1}{L_n} \sum_{r=L_n+\mathbb{T}_n^--1}^k \left[\left(\sum_{s=r-L_n+1}^r \hat{R}_{s,n} \right)^{\otimes 2} - \left(\sum_{s=r-L_n+1}^r R_{P,s,n} \right)^{\otimes 2} \right] \right\|_{\text{tr}} \wedge \epsilon \right) \\ &\stackrel{(2)}{\leq} \frac{1}{L_n} \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{E}_P \left(\left[\sum_{r \in \mathcal{T}_{n,L}} \left\| \left(\sum_{s=r-L_n+1}^r \hat{R}_{s,n} \right)^{\otimes 2} - \left(\sum_{s=r-L_n+1}^r R_{P,s,n} \right)^{\otimes 2} \right\|_{\text{tr}} \right] \wedge \epsilon \right) \\ &\stackrel{(3)}{\leq} \frac{2}{L_n} \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{E}_P \left(\left[\sum_{r \in \mathcal{T}_{n,L}} \left(\left\| \sum_{s=r-L_n+1}^r (\hat{R}_{s,n} - R_{P,s,n}) \right\|_2 \left\| \sum_{s=r-L_n+1}^r R_{P,s,n} \right\|_2 \right. \right. \right. \\ &\quad \left. \left. \left. + \left\| \sum_{s=r-L_n+1}^r (\hat{R}_{s,n} - R_{P,s,n}) \right\|_2^2 \right) \right] \wedge \epsilon \right), \end{aligned}$$

where (1) is from the definitions of $Q_{P,k,n}^R$, $\hat{Q}_{k,n}^R$, (2) is from the triangle inequality, and (3) is from the following outer product inequality for vectors $\hat{v}, v \in \mathbb{R}^d$

$$\begin{aligned} &\|\hat{v}\hat{v}^\top - vv^\top\|_{\text{tr}} \\ &\stackrel{(1)}{=} \|(\hat{v}-v)v^\top + v(\hat{v}-v)^\top + (\hat{v}-v)(\hat{v}-v)^\top\|_{\text{tr}} \\ &\stackrel{(2)}{\leq} 2\|(\hat{v}-v)v^\top\|_{\text{tr}} + \|(\hat{v}-v)(\hat{v}-v)^\top\|_{\text{tr}} \\ &\stackrel{(3)}{=} 2\|\hat{v}-v\|_2\|v\|_2 + \|\hat{v}-v\|_2^2, \end{aligned}$$

where (1) follows from adding and subtracting terms, (2) follows from the triangle inequality, and (3) follows by the properties of outer products and the definition of the trace norm. For each $r \in \mathcal{T}_{n,L}$, we have the following decomposition into the three bias terms from Step 1 by the triangle inequality

$$\begin{aligned} & \left\| \sum_{s=r-L_n+1}^r \left(\hat{\mathbf{R}}_{s,n} - \mathbf{R}_{P,s,n} \right) \right\|_2 \\ & \leq \left\| \sum_{s=r-L_n+1}^r \hat{\mathbf{w}}_{P,s,n}^{f,g}(\mathbf{Z}_{s,n}) \right\|_2 + \left\| \sum_{s=r-L_n+1}^r \hat{\mathbf{w}}_{P,s,n}^{g,\varepsilon}(\mathbf{Z}_{s,n}) \right\|_2 + \left\| \sum_{s=r-L_n+1}^r \hat{\mathbf{w}}_{P,s,n}^{f,\xi}(\mathbf{Z}_{s,n}) \right\|_2. \end{aligned}$$

Observe that for any $\delta > 0$ and any $r \in \mathcal{T}_{n,L}$, we have

$$\begin{aligned} & \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P \left(\left\| \sum_{s=r-L_n+1}^r \hat{\mathbf{w}}_{P,s,n}^{f,g}(\mathbf{Z}_{s,n}) \right\|_2 > \delta L_n^{\frac{1}{2}} \tau_n^7 D_n^{-2} \right) \\ & \leq \delta^{-1} L_n^{-\frac{1}{2}} \tau_n^{-7} D_n^2 \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{E}_P \left(\left\| \sum_{s=r-L_n+1}^r \hat{\mathbf{w}}_{P,s,n}^{f,g}(\mathbf{Z}_{s,n}) \right\|_2 \right) \\ & \leq \delta^{-1} L_n^{\frac{1}{2}} \tau_n^{-7} D_n^3 \sup_{P \in \mathcal{P}_{0,n}^*} \max_{(i,j,a,b) \in \mathcal{D}_n} \mathbb{E}_P(|\hat{w}_{P,t,n,i,a}^f|^2)^{\frac{1}{2}} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g|^2)^{\frac{1}{2}} \\ & = o(1), \end{aligned}$$

using the same arguments as Step 1.1 replacing $T_{n,L}$ with L_n , and noting that $D_n = O(T_n^{\frac{1}{6}})$ which corresponds to a lag-window size of $L_n = O(T_n^{\frac{1}{3-\delta'}})$ for any $\delta' > 0$. Next, for any $\delta > 0$ and any $r \in \mathcal{T}_{n,L}$, we have

$$\begin{aligned} & \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P \left(\left\| \sum_{s=r-L_n+1}^r \hat{\mathbf{w}}_{P,s,n}^{g,\varepsilon}(\mathbf{Z}_{s,n}) \right\|_2 > \delta L_n^{\frac{1}{2}} D_n^{-2} \tau_n^7 \right) \\ & = \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P \left(\left\| \sum_{s=r-L_n+1}^r \hat{\mathbf{w}}_{P,s,n}^{g,\varepsilon}(\mathbf{Z}_{s,n}) \right\|_2^2 > \delta^2 L_n D_n^{-4} \tau_n^{14} \right) \\ & \leq \delta^{-2} L_n^{-1} D_n^4 \tau_n^{-14} \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{E}_P \left(\left\| \sum_{s=r-L_n+1}^r \hat{\mathbf{w}}_{P,s,n}^{g,\varepsilon}(\mathbf{Z}_{s,n}) \right\|_2^2 \right) \\ & \leq \delta^{-2} D_n^5 \tau_n^{-14} \bar{K}^2 \sup_{P \in \mathcal{P}_{0,n}^*} \max_{t \in \mathcal{T}_{n,L}} \max_{j \in [d_Y]} \max_{b \in B_j} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g|^2) \\ & = o(1), \end{aligned}$$

for some $\bar{K} > 0$ using the same arguments as Step 1.2 replacing $T_{n,L}$ with L_n . The same arguments as Step 1.2 (i.e. exchanging g, ε with f, ξ) can be used to show that

$$\left\| \sum_{s=r-L_n+1}^r \hat{\mathbf{w}}_{P,s,n}^{f,\xi}(\mathbf{Z}_{s,n}) \right\|_2 = o_P(L_n^{\frac{1}{2}} D_n^{-2} \tau_n^7).$$

Hence, for any $r \in \mathcal{T}_{n,L}$ we have

$$\left\| \sum_{s=r-L_n+1}^r \left(\hat{\mathbf{R}}_{s,n} - \mathbf{R}_{P,s,n} \right) \right\|_2 = o_P(L_n^{\frac{1}{2}} D_n^{-2} \tau_n^7).$$

By Lemma B.2, we have for all $r \in \mathcal{T}_{n,L}$ that

$$\begin{aligned}
& \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P \left(\left\| \sum_{s=r-L_n+1}^r \mathbf{R}_{P,s,n} \right\|_2 > L_n^{\frac{1}{2}} D_n^{\frac{1}{2}} \epsilon \right) \\
&= \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P \left(\left\| \sum_{s=r-L_n+1}^r \mathbf{R}_{P,s,n} \right\|_2^2 > L_n D_n \epsilon^2 \right) \\
&\leq L_n^{-1} D_n^{-1} \epsilon^{-2} \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{E}_P \left(\left\| \sum_{s=r-L_n+1}^r \mathbf{R}_{P,s,n} \right\|_2^2 \right) \\
&\leq L_n^{-1} D_n^{-1} \epsilon^{-2} (2 L_n^{\frac{1}{2}} D_n^{\frac{1}{2}} \bar{\Theta}^R K \sum_{h=1}^{\infty} h^{-\bar{\beta}^R})^2 \\
&= O(1),
\end{aligned}$$

where $\sum_{h=1}^{\infty} h^{-\bar{\beta}^R} < \infty$ since $\bar{\beta}^R > 1$ by Assumption 3.5. Putting it all together, by Markov's inequality, bounded convergence (Lemma B.7), and noting that the previous statements hold for all times in $\mathcal{T}_{n,L}$, we have

$$\max_{k \in \mathcal{T}_{n,L}} \|\hat{Q}_{k,n}^{\mathbf{R}} - Q_{P,k,n}^{\mathbf{R}}\|_{\text{tr}} = O_{\mathcal{P}}(r_{n,2}^{\delta}),$$

where

$$r_{n,2}^{\delta} = T_{n,L} \tau_n^7 D_n^{-\frac{3}{2}} + T_{n,L} D_n^{-4} \tau_n^{14}.$$

Next, recall the offset condition

$$\nu_n \gg \log(T_{n,L}) D_n (\tau_n^{-2} (T_{n,L}^{-\frac{1}{2}} D_n^{-\frac{1}{2}} \hat{\delta}_{P,n}^{\frac{1}{2}} + T_{n,L}^{-1})),$$

where

$$\hat{\delta}_{P,n} = O_{\mathcal{P}}(r_{n,1}^{\delta} + r_{n,2}^{\delta}),$$

and

$$\begin{aligned}
r_{n,1}^{\delta} &= D_n (\bar{\Gamma}_n^R L_n^{\frac{1}{2}} + T_{n,L}^{\frac{1}{2}} D_n^{\frac{1}{2}} L_n^{\frac{1}{2}} + T_{n,L} L_n^{-1} + T_{n,L} L_n^{2-\bar{\beta}^R}), \\
r_{n,2}^{\delta} &= T_{n,L} \tau_n^7 D_n^{-\frac{3}{2}} + T_{n,L} D_n^{-4} \tau_n^{14}.
\end{aligned}$$

Observe that

$$\begin{aligned}
& T_{n,L}^{-\frac{1}{2}} D_n^{-\frac{1}{2}} (r_{n,1}^{\delta})^{\frac{1}{2}} + T_{n,L}^{-1} \\
&\leq T_{n,L}^{-\frac{1}{2}} D_n^{-\frac{1}{2}} (D_n^{\frac{1}{2}} ((\bar{\Gamma}_n^R)^{\frac{1}{2}} L_n^{\frac{1}{4}} + T_{n,L}^{\frac{1}{4}} D_n^{\frac{1}{4}} L_n^{\frac{1}{4}} + T_{n,L}^{\frac{1}{2}} L_n^{-\frac{1}{2}} + T_{n,L}^{\frac{1}{2}} L_n^{1-\frac{\bar{\beta}^R}{2}})) + T_{n,L}^{-1} \\
&= T_{n,L}^{-\frac{1}{2}} (\bar{\Gamma}_n^R)^{\frac{1}{2}} L_n^{\frac{1}{4}} + T_{n,L}^{-\frac{1}{4}} D_n^{\frac{1}{4}} L_n^{\frac{1}{4}} + L_n^{-\frac{1}{2}} + L_n^{1-\frac{\bar{\beta}^R}{2}} + T_{n,L}^{-1} \\
&= \varphi_{n,1},
\end{aligned}$$

which comes from the covariance estimation error. Also, we have

$$\begin{aligned}
& T_{n,L}^{-\frac{1}{2}} D_n^{-\frac{1}{2}} (r_{n,2}^{\delta})^{\frac{1}{2}} \\
&\leq T_{n,L}^{-\frac{1}{2}} D_n^{-\frac{1}{2}} (T_{n,L}^{\frac{1}{2}} \tau_n^{\frac{7}{2}} D_n^{-\frac{3}{4}} + T_{n,L}^{\frac{1}{2}} \tau_n^7 D_n^{-2}) \\
&= \tau_n^{\frac{7}{2}} D_n^{-\frac{5}{4}} + \tau_n^7 D_n^{-\frac{5}{2}} \\
&= \varphi_{n,2},
\end{aligned}$$

which comes from the prediction errors since we use the residual products instead of the error products. The assumption on the offset condition (13) implies that

$$\nu_n \gg \log(T_{n,L}) D_n (\tau_n^{-2} (\varphi_{n,1} + \varphi_{n,2})),$$

and therefore

$$\sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(\hat{\Delta}_{P,n} \geq \frac{1}{16} T_{n,L} \nu_n \tau_n^2) \rightarrow 0.$$

Combining the results from Step 1, Step 2, and Step 3, we obtain the final result

$$\limsup_{n \rightarrow \infty} \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{P}_P(S_{n,p}(\hat{\mathbf{R}}_n) > \hat{q}_{1-\alpha+\nu_n} + \tau_n) \leq \alpha.$$

□

C.2 Proof of Theorem 4.1

It suffices to establish the following two points. First, that the assumptions of Theorem 4.1 imply those of Theorem 3.1. Second, that the sieve time-varying regression estimators, under the setup of Theorem 4.1, satisfy the convergence rate requirements of Theorem 3.1.

By using the following notation, we see that Assumption 4.1 implies Assumption 3.1 and Assumption 4.4 implies Assumption 3.4. Note that the causal representations for the observed processes and error processes from Assumptions 4.1 and 4.4 are defined for all rescaled times \mathcal{U}_n , and therefore for all $\{t/n\}_{t \in \mathcal{T}_n} \subset \mathcal{U}_n$ in particular. For a generic high-dimensional locally stationary observed process $W \in \{X, Y, Z\}$ and any time t , sample size n , dimension l , and time-offset d , we write

$$G_{t,n}^W(\cdot) = \tilde{G}_n^W(t/n, \cdot), \quad G_{t,n,l}^W(\cdot) = \tilde{G}_{n,l}^W(t/n, \cdot), \quad G_{t,n,l,d}^W(\cdot) = \tilde{G}_{n,l,d}^W(t/n, \cdot),$$

to respectively denote the causal representations of all dimensions of the process W , dimension l of the process W , and dimension l of the process W with time-offset d . For a generic high-dimensional locally stationary error process $e \in \{\varepsilon, \xi\}$ and any distribution P , time t , sample size n , dimension l , and time-offset d , we denote

$$G_{P,t,n}^e(\cdot) = \tilde{G}_{P,n}^e(t/n, \cdot), \quad G_{P,t,n,l}^e(\cdot) = \tilde{G}_{P,n,l}^e(t/n, \cdot), \quad G_{P,t,n,l,d}^e(\cdot) = \tilde{G}_{P,n,l,d}^e(t/n, \cdot),$$

to respectively denote the causal representations of all dimensions of the error process e , dimension l of the error process e , and dimension l of the error process e with time-offset d . The causal representations of the error products are defined similarly. Using this notation, we see that Assumption 4.5 implies Assumption 3.5 and Assumption 4.6 implies Assumption 3.6. Specifically, Assumption 3.6 is satisfied with $\bar{\Gamma}_n^R = D_n^{\frac{1}{2}}$ by using linearity of expectation and directly applying the stochastic Lipschitz condition for the product of errors from the discussion below Assumption 4.6 to each term in the sum.

It remains to show that Assumptions 3.2 and 3.3 are implied. To see this, let us consider the following notation. For any distribution P , time t , sample size n , dimensions i, j , and time-offsets a, b , we write

$$\begin{aligned} f_{P,t,n,i,a}(\cdot) &= f_{P,n,i,a}(t/n, \cdot), \quad \hat{f}_{t,n,i,a}(\cdot) = \hat{f}_{t,n,i,a}(t/n, \cdot), \\ g_{P,t,n,j,b}(\cdot) &= g_{P,n,j,b}(t/n, \cdot), \quad \hat{g}_{t,n,j,b}(\cdot) = \hat{g}_{t,n,j,b}(t/n, \cdot), \end{aligned}$$

to denote the time-varying regression functions and the corresponding sieve estimators from Subsection 4.3 using the notation of Subsection 2.3. For all times $t \in \mathcal{T}_n$, the algorithms used to construct the sieve estimators from Subsection 4.3 for rescaled time $t/n \in \mathcal{U}_n$ are Borel measurable functions of the datasets $\mathfrak{D}_{t,n,i,a}^{\hat{f}}$ and $\mathfrak{D}_{t,n,j,b}^{\hat{g}}$. The measurability of the causal mechanisms of the observed processes from Assumption 4.1 ensures that these sieve estimators have the causal representations $G_{t,n,i,a}^{\mathcal{A}^{\hat{f}}}(\mathcal{H}_{t,a}^{\mathfrak{D}^{\hat{f}}})$ and $G_{t,n,j,b}^{\mathcal{A}^{\hat{g}}}(\mathcal{H}_{t,b}^{\mathfrak{D}^{\hat{g}}})$ from Assumption 3.2.

Further, note that the sieve estimators are Borel measurable functions from \mathbb{R}^{dz} to \mathbb{R} . The measurability of the causal mechanisms of the covariate processes from Assumption 4.1 ensures that the sieve estimator's predictions have the causal representations $G_{t,n,i,a}^{\hat{f}}(\mathcal{H}_{t,a}^{\hat{f}})$ and $G_{t,n,j,b}^{\hat{g}}(\mathcal{H}_{t,b}^{\hat{g}})$ from Assumption 3.3. Similarly, note that the Borel measurability of the conditional expectations $f_{P,t,n,i,a}$ and $g_{P,t,n,j,b}$ ensures that the sieve estimator's prediction errors are Borel measurable functions from \mathbb{R}^{dz} to \mathbb{R} . Again, by the measurability of the causal mechanisms of the covariate processes from Assumption 4.1, the prediction errors are ensured to have the causal representations $G_{P,t,n,i,a}^{\hat{w}^{\hat{f}}}(\mathcal{H}_{t,a}^{\hat{f}})$

and $G_{P,t,n,j,b}^{\hat{w}^g}(\mathcal{H}_{t,b}^{\hat{g}})$ from Assumption 3.3. In view of the boundedness of the sieve estimator's predictions by construction, the regularity conditions for the time-varying partial response functions from Assumption 4.2, and the additive form of the time-varying regression functions from Assumption 4.2, there exists some $q \geq 2$ such that for all $n \in \mathbb{N}$, $t \in \mathcal{T}_n$, and $(i, j, a, b) \in \mathcal{D}_n$, the prediction errors satisfy

$$\sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{E}_P(|\hat{w}_{P,t,n,i,a}^f|^q) < \infty, \quad \sup_{P \in \mathcal{P}_{0,n}^*} \mathbb{E}_P(|\hat{w}_{P,t,n,j,b}^g|^q) < \infty.$$

Hence, the sieve estimator's predictions and prediction errors meet all the conditions required by Assumption 3.3.

The distribution-uniform assumptions of Theorem 4.1 imply that the distribution-pointwise assumptions of Theorem 3.2 in Ding and Zhou [48] hold for each distribution in the collection. Specifically, for each $n \in \mathbb{N}$ and $P \in \mathcal{P}_{0,n}^*$, Assumption 4.2 implies the additive form of the time-varying regression functions in [48], Assumptions 4.1, 4.4, 4.6 imply Assumption 2.1 in [48], Assumption 4.5 implies Assumption 2.2 in [48], Assumption 4.3 implies Assumption 3.1 in [48], and Assumption 4.7 implies Assumption 3.2 in [48]. Next, we consider the additional regularity condition required by Theorem 3.2 in Ding and Zhou [48] involving the rate of decay in temporal dependence and the rate of growth of the largest sup-norm of the basis functions for time.

Recall the numbers of observations $T_{t,n,i,a}^{\hat{f}}, T_{t,n,j,b}^{\hat{g}}$ in the datasets $\mathfrak{D}_{t,n,i,a}^{\hat{f}}, \mathfrak{D}_{t,n,j,b}^{\hat{g}}$ used to construct the sieve estimators $\hat{f}_{t,n,i,a}(t/n, \cdot), \hat{g}_{t,n,j,b}(t/n, \cdot)$ of the time-varying regression functions at rescaled time $t/n \in \mathcal{U}_n$. Also, recall the numbers of basis functions \tilde{c}_n, \tilde{d}_n from Subsection 4.3. As previously noted in Subsection 4.3, we simplified the notation for the numbers of basis functions $\{\phi_{\ell_1}(u)\}, \{\varphi_{\ell_2}(z)\}$ for the estimators $\hat{f}_{t,n,i,a,k,c}(t/n, \cdot)$ and $\hat{g}_{t,n,j,b,k,c}(t/n, \cdot)$ of the time-varying partial response functions at rescaled time $t/n \in \mathcal{U}_n$ from $\tilde{c}_{t,n,i,a,k,c}^{\hat{f}}, \tilde{d}_{t,n,i,a,k,c}^{\hat{f}}$ and $\tilde{c}_{t,n,j,b,k,c}^{\hat{g}}, \tilde{d}_{t,n,j,b,k,c}^{\hat{g}}$ to \tilde{c}_n, \tilde{d}_n . We will now require the full notation for the numbers of basis functions.

For the convergence rate guarantees from Theorem 3.2 in Ding and Zhou [48] to be applicable in our setting, we must have

$$\begin{aligned} \tilde{c}_{t,n,i,a,k,c}^{\hat{f}} \tilde{d}_{t,n,i,a,k,c}^{\hat{f}} & \left(\frac{1}{\sqrt{T_{t,n,i,a}^{\hat{f}}}} + \frac{(T_{t,n,i,a}^{\hat{f}})^{\frac{\min(\bar{\beta}, \bar{\beta}^\infty)}{2}+1}}{T_{t,n,i,a}^{\hat{f}}} \right) \sup_{\ell_1 \in [\tilde{c}_{t,n,i,a,k,c}^{\hat{f}}]} \sup_{u \in \mathcal{U}_n} |\phi_{\ell_1}(u)|^2 = o(1), \\ \tilde{c}_{t,n,j,b,k,c}^{\hat{g}} \tilde{d}_{t,n,j,b,k,c}^{\hat{g}} & \left(\frac{1}{\sqrt{T_{t,n,j,b}^{\hat{g}}}} + \frac{(T_{t,n,j,b}^{\hat{g}})^{\frac{\min(\bar{\beta}, \bar{\beta}^\infty)}{2}+1}}{T_{t,n,j,b}^{\hat{g}}} \right) \sup_{\ell_1 \in [\tilde{c}_{t,n,j,b,k,c}^{\hat{g}}]} \sup_{u \in \mathcal{U}_n} |\phi_{\ell_1}(u)|^2 = o(1), \end{aligned}$$

for each time $t \in \mathcal{T}_n$ and combination of dimensions $i \in [d_X], j \in [d_Y], k \in [d_Z]$ and time-offsets $a \in A_i, b \in B_j, c \in C_k$. This condition is satisfied for the following reasons. First, we have

$$\sup_{\ell_1 \in [\tilde{c}_{t,n,i,a,k,c}^{\hat{f}}]} \sup_{u \in \mathcal{U}_n} |\phi_{\ell_1}(u)|^2 \lesssim (\tilde{c}_{t,n,i,a,k,c}^{\hat{f}})^2, \quad \sup_{\ell_1 \in [\tilde{c}_{t,n,j,b,k,c}^{\hat{g}}]} \sup_{u \in \mathcal{U}_n} |\phi_{\ell_1}(u)|^2 \lesssim (\tilde{c}_{t,n,j,b,k,c}^{\hat{g}})^2,$$

because the basis functions are chosen to be mapped Legendre polynomials; see Appendix C in Ding and Zhou [48] and Section 3 in Belloni et al. [14]. For more information about sieve estimators and other basis functions, see [113, 74, 29, 46, 47]. Second, because we have chosen the numbers of basis functions to be $O(\log(T_n))$ in the setup of Theorem 4.1. Third, because the constants $\bar{\beta}, \bar{\beta}^\infty$ from Assumption 4.5 are both greater than 2. Fourth, because $T_{t,n,i,a}^{\hat{f}} = o(T_n)$ and $T_{t,n,j,b}^{\hat{g}} = o(T_n)$ regardless of whether the sieve estimators are fit once based on all the data or sequentially as in Remark 4.1. To be clear, this is due to the infill asymptotic framework of locally stationary processes, so that more and more observations are available for each local structure as n grows.

Therefore, the main inequality in the proof of Theorem 3.2 in [48] holds for each $P \in \mathcal{P}_{0,n}^*$ and $n \in \mathbb{N}$ because all of the theorem's assumptions are satisfied under the stronger assumptions of Theorem 4.1. Moreover, for each $n \in \mathbb{N}$, the supremum over $P \in \mathcal{P}_{0,n}^*$ of the final upper bound for the main inequality in the proof of Theorem 3.2 in [48] is finite under the distribution-uniform assumptions of Theorem 4.1. Thus, by basic properties of the supremum, the main inequality in the proof of Theorem 3.2 in [48] holds with a supremum over $P \in \mathcal{P}_{0,n}^*$ for each $n \in \mathbb{N}$. In view of the notational changes described in Remark 3.2 in [48], the same steps in the proof of Theorem 3.2 in [48] imply that this distribution-uniform inequality also holds in the general regression setting with time-offsets. We do not repeat the

proof of Theorem 3.2 in [48] here, as the only changes are in the notation. Putting it all together, the prediction errors of the sieve estimators with the setup of Theorem 4.1 satisfy

$$\sup_{P \in \mathcal{P}_{0,n}^*} \max_{i \in [d_X], a \in A_i} \max_{t \in \mathcal{T}_n} \mathbb{E}_P \left(\left| \hat{w}_{P,t,n,i,a}^f \right|^2 \right)^{\frac{1}{2}} = o(T_n^{-\frac{1}{2+\delta}} \log^3(T_n)),$$

$$\sup_{P \in \mathcal{P}_{0,n}^*} \max_{j \in [d_Y], b \in B_j} \max_{t \in \mathcal{T}_n} \mathbb{E}_P \left(\left| \hat{w}_{P,t,n,j,b}^g \right|^2 \right)^{\frac{1}{2}} = o(T_n^{-\frac{1}{2+\delta}} \log^3(T_n)),$$

for any $\delta > 0$. Since $D_n = O(T_n^{\frac{1}{6}})$ and $\tau_n = o(\log^{-(1+\delta')}(T_n))$ for some $\delta' > 0$, the convergence rates required by Theorem 3.1 are achieved by the sieve estimators with the setup of Theorem 4.1.

□