

# A better understanding of Granger causality analysis: A big data environment\*

Xiaojun Song<sup>†</sup>

Abderrahim Taamouti<sup>‡</sup>

Peking University

Durham University Business School

February 15, 2017

## ABSTRACT

We provide a better understanding of the causal structure in a multivariate time series by introducing a novel statistical procedure for testing indirect and spurious causal effects. In practice, detecting these effects is a complicated task, since the auxiliary variables that transmit/induce indirect/spurious causality are often unknown. We propose an efficient statistical procedure to test for the presence of indirect/spurious causality based on big data analysis. We suggest an identification procedure to find the variables that transmit/induce the indirect/spurious causality. Finally, we provide an application where 135 economic variables were used to study a possible indirect causality from money/credit to income.

**Keywords:** Indirect causality, spurious causality, big data analysis, auxiliary variable(s), asymptotic theory, Monte Carlo simulations, money, credit and income.

**Journal of Economic Literature classification:** C12; C32; C38; C53; E60.

---

\*We would like to thank very much Raffaella Giacomini for her discussion which helped us write this paper.

<sup>†</sup>Department of Business Statistics and Econometrics, Guanghua School of Management and Center for Statistical Science, Peking University, Beijing, 100871, China. E-mail: sxj@gsm.pku.edu.cn. Financial support from the National Natural Science Foundation of China (Grant No. 71532001) is acknowledged.

<sup>‡</sup>Department of Economics and Finance, Durham University Business School. Address: Mill Hill Lane, Durham, DH1 3LB, UK. TEL: +44-1913345423. E-mail: abderrahim.taamouti@durham.ac.uk.

# 1 Introduction

The concept of causality introduced by Wiener (1956) and Granger (1969) constitutes a basic notion for analyzing dynamic relationships between time series. In studying Wiener-Granger causality, predictability is the central issue, hence its importance to economists and policymakers. In practice, Granger-causality is often investigated for bivariate processes. However, different conclusions may be reached when more than two variables are considered. If more than two variables are present, non-causality conditions become more complicated; see Lütkepohl (1993) and Dufour and Renault (1998). In other words, even if a variable is Granger-causal in a bivariate model, it may not be Granger-causal in a larger model involving more variables. In this case, we talk about an *indirect causality* transmitted through a third variable(s); hereafter auxiliary variable(s). For instance, there may be a variable that drives both variables in the bivariate process, thus when this variable is included into the model, a bivariate causal structure may disappear. In turn, it is also possible that a variable is non-causal for another one in a bivariate model and becomes causal if the information set is extended to include other variables as well. The latter situation corresponds to what is known as a *spurious causality*. Ignoring these causal effects can lead to wrong economic analysis, and consequently to inaccurate policy decisions. In this paper, we introduce a novel statistical procedure that allows us to test for indirect and spurious causal effects.

Wiener-Granger's concept of causality is defined in terms of predictability of a variable  $Y$  from its own past and the past of another variable  $X$ . The original definition of Granger (1969) implicitly assumes that all the relevant information is available and used for the causality analysis. However, in practice only a very limited information is considered and the omission of key variables (auxiliary variables) could lead to a spurious causality and might not help detect a possible indirect causality between the variables of interest. The relevance of the information set for Granger causality analysis was first pointed out by Hsiao (1982), who formally introduced the concept of indirect/spurious causality in a trivariate model. Hsiao (1982) provides a basic framework to explain the causal relationships in a multivariate time series model based on Wiener-Granger notion of causality. He focuses on establishing a Granger causal ordering of the events and on the reconciliation of the disparity between the results obtained from the bivariate and multivariate analysis. He generalizes the Granger's concept of causality to make some provision for spurious/indirect causality which may arise in multivariate analysis. In particular, he shows that a certain type of spurious causality vanishes when the information set is reduced. This observation leads to a strengthened definition of

(direct) causality by requiring an improvement in prediction irrespective of the used information set. Finally, Hsiao (1982) characterizes the indirect/spurious causality in the context of VAR models and discusses how to test these causal effects in the presence of *known* auxiliary variables.

The main issue of Hsiao (1982)'s framework is that the auxiliary variables that transmit/induce the indirect/spurious causality are implicitly assumed to be *known*. However, in practice these variables are unknown, except in the presence of an economic theory that explicitly specifies the auxiliary variables, which complicates very much the task of testing for the presence of an indirect/spurious causality. The availability of hundreds of economic variables makes this task even harder as it is generally infeasible to find the appropriate auxiliary variable(s) among all the available ones. In addition, including hundreds of variables and their lags in a regression equation is technically difficult.

In this paper, we introduce new statistical procedures to test for the presence of indirect/spurious causality using big data analysis. To overcome the problem of unknown relevant auxiliary variables, a diffusion index, extracted using principal component analysis, is included in the regression equation to represent all the variables that are available to practitioners. We derive the asymptotic distributions of the tests in the presence of the estimated index. Furthermore, we conduct a Monte Carlo simulation to evaluate the performance of the proposed statistical procedures. The results show that our procedures are efficient for detecting indirect/spurious causality.

Unfortunately, the above statistical procedures only test for the presence/absence of an indirect/spurious causality and cannot inform us about the variables of the big data that are responsible for the transmission/induction of this indirect/spurious causality. Another contribution of this paper is we provide an identification procedure which helps us identify the variables in the big data that transmit/induce the indirect/spurious causality.

Finally, to show the practical relevance of the proposed tests, we use 135 economic variables to examine the causality from money/credit to income. In particular, we test whether or not there is an *indirect* causality from monetary policy/credit to income. Thereafter, if this indirect causality exists, then we use the identification procedure discussed above to identify the auxiliary variable(s) that are responsible for the transmission of this indirect causality. Our empirical results show that there is an indirect causality from credit to income, but not from money to income. In addition, the identification procedure indicates that this indirect causality is mainly transmitted through short and long-term interest rates. Hence, interest rates are responsible for the indirect causality from credit to income.

The plan of the paper is as follows. Section 2 presents the general theoretical framework which

underlies the definition of indirect/spurious causality. Section 3 provides some motivations for deriving statistical procedures that help detect indirect/spurious causality. In Section 4, we define the regression models and hypotheses tests that we consider to test for indirect/spurious causality. In Section 5, we provide the asymptotic distributions of the tests. These distributions are derived based on the asymptotic theory from the factor analysis. In Section 6, we propose a statistical procedure that allows us to identify the auxiliary variables that transmit/induce the indirect/spurious causality. In Section 7, we run a Monte Carlo simulation to investigate the finite sample properties of the tests of indirect/spurious causality. Section 8 is devoted to an empirical application. The conclusion is given in Section 9. The proofs of the theoretical results are presented in Appendix A. Finally, a description of all variables used in Section 8 can be found in a separate companion Appendix, which is available online [see Song and Taamouti (2017)].

## 2 Framework

We consider three stochastic processes  $\{X_t : t \in \mathbb{Z}\}$ ,  $\{Y_t : t \in \mathbb{Z}\}$ , and  $\{Z_t : t \in \mathbb{Z}\}$ . For simplicity of exposition, we assume that these processes are univariate. We denote  $I_X(t) = \{X(s) : s \leq t\}$ ,  $I_Y(t) = \{Y(s) : s \leq t\}$  and  $I_Z(t) = \{Z(s) : s \leq t\}$  the information sets which contain all the past and present values of  $X$ ,  $Y$ , and  $Z$  until time  $t$ , respectively. We denote  $I(t)$  the information set that contains  $I_X(t)$ ,  $I_Y(t)$  and  $I_Z(t)$ .  $I(t) - A_t$ , with  $A_t = I_X(t), I_Y(t), I_Z(t)$ , contains all the elements of  $I(t)$  except those of  $A_t$ . The notion of non-causality considered here is defined in terms of orthogonality conditions between subspaces of a Hilbert space of random variables with finite second moments. We denote  $L^2 \equiv L^2(\Omega, \mathcal{A}, Q)$  the Hilbert space of real random variables defined on a common probability space  $(\Omega, \mathcal{A}, Q)$ , with covariance as inner product.

For any information set  $B_t$  [some Hilbert subspace of  $L^2$ ], we denote  $P[X_{t+1} | B_t]$  the best linear forecast of  $X_{t+1}$  based on the information set  $B_t$ , the corresponding prediction error is

$$u(X_{t+1} | B_t) = X_{t+1} - P[X_{t+1} | B_t],$$

and  $\sigma^2(X_{t+1} | B_t)$  is the variance of the prediction error.  $P[X_{t+1} | B_t]$  is the orthogonal projection of  $X_{t+1}$  on the subspace  $B_t$ .

We now remind the reader of the following definitions of *indirect* causality and *spurious* causality from Hsiao (1982). In the following, the random variable  $Z$  is used as a *known* auxiliary variable. However, in the next sections, when we describe the statistical procedures for testing indirect/spurious causality,  $Z$  will be treated as an *unknown* auxiliary variable.

**Definition 1** [*Indirect Causality*]:  $Y$  is an indirect cause of  $X$ , denoted  $Y \overset{ind}{\mapsto} X | I_X(t), I_Z(t)$ , iff

(i)  $Y$  Granger causes  $X$  with respect to the information set  $I_X(t)$  :

$$P[X_{t+1} | I_X(t)] \neq P[X_{t+1} | I(t) - I_Z(t)], \text{ for some } t > w,$$

(ii)  $Y$  does not Granger cause  $X$  with respect to the information set  $I(t) - I_Y(t)$  :

$$P[X_{t+1} | I(t) - I_Y(t)] = P[X_{t+1} | I(t)], \forall t > w,$$

(iii): (a)  $Y$  Granger causes  $Z$  and (b)  $Z$  Granger causes  $X$  with respect to the information sets  $I(t) - I_Y(t)$  and  $I(t) - I_Z(t)$ , respectively:

$$P[Z_{t+1} | I(t) - I_Y(t)] \neq P[Z_{t+1} | I(t)], \text{ for some } t > w,$$

$$P[X_{t+1} | I(t) - I_Z(t)] \neq P[X_{t+1} | I(t)], \text{ for some } t > w,$$

where  $w$  is a “starting point” which is typically equal to a finite initial date [such as  $w = -1, 0$ , or  $1$ ] or to  $-\infty$ ; in the latter case  $I(t)$  is defined for all  $t \in \mathbb{Z}$ .

Thus, the conditions [(i), (ii), (iii)] of Definition 1 must be satisfied in order to have an indirect causality from  $Y$  to  $X$  in the presence of an auxiliary variable  $Z$ . Similar conditions can be obtained for an indirect causality from  $X$  to  $Y$ . We now provide the necessary conditions for a spurious causality from  $Y$  to  $X$ . We distinguish between two types of spurious causality.

**Definition 2** [*Spurious Causality*]:

1.  $Y$  is a spurious cause of type I for  $X$  if

1.(i)  $Y$  Granger causes  $X$  with respect to the information set  $I(t) - I_Y(t)$  :

$$P[X_{t+1} | I(t) - I_Y(t)] \neq P[X_{t+1} | I(t)], \text{ for some } t > w,$$

1.(ii)  $Y$  does not Granger cause  $X$  with respect to the information set  $I_X(t)$  :

$$P[X_{t+1} | I_X(t)] = P[X_{t+1} | I(t) - I_Z(t)], \forall t > w,$$

1.(iii): (a)  $Y$  Granger causes  $Z$  and (b)  $Z$  Granger causes  $X$ , both with respect to the information sets  $I(t) - I_Y(t)$  and  $I(t) - I_Z(t)$ , respectively,

$$P[Z_{t+1} | I(t) - I_Y(t)] \neq P[Z_{t+1} | I(t)], \text{ for some } t > w,$$

$$P[X_{t+1} | I(t) - I_Z(t)] \neq P[X_{t+1} | I(t)], \text{ for some } t > w.$$

2.  $Y$  is a spurious cause of type II for  $X$  if

2. **(i)**  $Y$  Granger causes  $X$  with respect to the information set  $I_X(t)$  :

$$P[X_{t+1} | I_X(t)] \neq P[X_{t+1} | I(t) - I_Z(t)], \text{ for some } t > w,$$

2. **(ii)**  $Y$  does not Granger cause  $X$  with respect to the information set  $I(t) - I_Y(t)$  :

$$P[X_{t+1} | I(t) - I_Y(t)] = P[X_{t+1} | I(t)], \forall t > w,$$

2. **(iii): (a)**  $Z$  Granger causes  $Y$  and **(b)**  $Z$  Granger causes  $X$ , both with respect to the information set  $I(t) - I_Z(t)$  :

$$P[Y_{t+1} | I(t) - I_Z(t)] \neq P[Y_{t+1} | I(t)], \text{ for some } t > w,$$

$$P[X_{t+1} | I(t) - I_Z(t)] \neq P[X_{t+1} | I(t)], \text{ for some } t > w.$$

Definition 2 shows that there are three conditions to satisfy for each type [type I and type II] of spurious causality from  $Y$  to  $X$ . Similar conditions can be obtained for the spurious causality from  $X$  to  $Y$ .

The previous definitions will be used to construct statistical procedures that test for the presence of indirect/spurious causality. We next assume that only the variables of interest  $X$  and  $Y$  are observable. Thus, the auxiliary variable  $Z$  will be treated as an unknown variable, and it will be extracted using factor analysis.

### 3 Motivation

Unfortunately, most empirical studies on Granger causality analysis ignore indirect and spurious causal effects. This might be explained by the lack of statistical procedures that detect these effects. Up to now, the detection of indirect/spurious causality depends on the knowledge of relevant auxiliary variables, which can happen only in rare cases such as the existence of an economic theory that identifies these variables. The following examples illustrate situations where indirect/spurious causality happens. We start with an example on indirect causality.

**Example 1 [Indirect Causality]** *This example illustrates an indirect causality between two variables of interest  $X$  and  $Y$  transmitted through a third auxiliary variable  $Z$ . We consider the following regressions with different information sets. In the absence of  $Z$ , we assume that*

$$X_t = \mu_X^{(1)} + \alpha_X^{(1)} X_{t-1} + \alpha_Y^{(1)} Y_{t-1} + \varepsilon_{X,t}^{(1)}, \text{ where } \varepsilon_{X,t}^{(1)} \sim N(0, 1). \quad (1)$$

In the presence of  $Z$ , we assume

$$X_t = \mu_X^{(2)} + \alpha_X^{(2)} X_{t-1} + \alpha_Y^{(2)} Y_{t-2} + \alpha_Z^{(2)} Z_{t-1} + \varepsilon_{X,t}^{(2)}, \text{ with } \alpha_Z^{(2)} \neq 0 \text{ and } \varepsilon_{X,t}^{(2)} \sim N(0, 1). \quad (2)$$

Now, if we assume that the auxiliary variable  $Z$  is generated by the process

$$Z_t = \mu_Z + \phi_Y Y_{t-1} + \phi_Z Z_{t-1} + \varepsilon_{Z,t}, \text{ with } \phi_Y = -\frac{\alpha_Y^{(2)}}{\alpha_Z^{(2)}} \text{ and where } \varepsilon_{Z,t} \sim N(0, 1), \quad (3)$$

then we have an **indirect causality** from  $Y$  to  $X$  which is transmitted by  $Z$ . We can easily check that all the conditions of Definition 1 are satisfied. In particular  $Y$  causes  $X$  in the absence of  $Z$  [see Equation (1)], but not in the presence of  $Z$ . The latter can be seen by plugging Equation (3) into Equation (2) to obtain

$$X_t = \mu_X^{(2)} + \alpha_X^{(2)} X_{t-1} + \alpha_Y^{(2)} Y_{t-2} + \alpha_Z^{(2)} [\mu_Z + \phi_Y Y_{t-2} + \phi_Z Z_{t-2} + \varepsilon_{Z,t-1}] + \varepsilon_{X,t}^{(2)}. \quad (4)$$

Since  $\phi_Y = -\frac{\alpha_Y^{(2)}}{\alpha_Z^{(2)}}$ , the term  $\alpha_Y^{(2)} Y_{t-2} + \alpha_Z^{(2)} \phi_Y Y_{t-2}$  in Equation (4) is equal to zero and

$$X_t = \left( \mu_X^{(2)} + \alpha_Z^{(2)} \mu_Z \right) + \alpha_X^{(2)} X_{t-1} + \alpha_Z^{(2)} \phi_Z Z_{t-2} + \alpha_Z^{(2)} \varepsilon_{Z,t-1} + \varepsilon_{X,t}^{(2)}.$$

Hence,  $Y$  does not cause  $X$  in the presence of  $Z$ .

**Example 2 [Spurious Causality]** This example illustrates a spurious causality of type I from  $Y$  to  $X$  induced by an auxiliary variable  $Z$ . We consider the following regression models with different information sets. In the absence of  $Z$ , we assume that

$$X_t = \mu_X^{(1)} + \phi_X^{(1)} X_{t-1} + \varepsilon_{X,t}^{(1)}, \text{ where } \varepsilon_{X,t}^{(1)} \sim N(0, 1). \quad (5)$$

In the presence of  $Z$ , we assume

$$X_t = \mu_X^{(2)} + \phi_X^{(2)} X_{t-1} + \phi_Z^{(2)} Z_{t-1} + \varepsilon_{X,t}^{(2)}, \text{ with } \phi_Z^{(2)} \neq 0 \text{ and } \varepsilon_{X,t}^{(2)} \sim N(0, 1). \quad (6)$$

Now, if we assume that the auxiliary variable  $Z$  is generated by the process

$$Z_t = \mu_Z + \phi_Z^{(3)} Z_{t-1} + \phi_Y^{(3)} Y_{t-1} + \phi_X^{(3)} X_{t-1} + \varepsilon_{Z,t}, \text{ where } \varepsilon_{Z,t} \sim N(0, 1), \quad (7)$$

then we have a **spurious causality of Type I** from  $Y$  to  $X$  induced by  $Z$ . We can easily check that all the conditions of Definition 2 are satisfied. In particular,  $Y$  does not cause  $X$  in the absence of  $Z$  [see Equation (5)], but it does in the presence of  $Z$ . The latter can be seen by plugging Equation (7) into Equation (6) to obtain:

$$X_t = \mu_X^{(2)} + \phi_Z^{(2)} \mu_Z + \phi_X^{(2)} X_{t-1} + \phi_Z^{(2)} \phi_X^{(3)} X_{t-2} + \phi_Z^{(2)} \phi_Z^{(3)} Z_{t-2} + \phi_Z^{(2)} \phi_Y^{(3)} Y_{t-2} + \phi_Z^{(2)} \varepsilon_{Z,t-1} + \varepsilon_{X,t}^{(2)}. \quad (8)$$

Hence, for  $\phi_Y^{(3)} \neq 0$  (in addition to  $\phi_Z^{(2)} \neq 0$ ), from Equation (8) we see that  $Y$  spuriously (Type I) causes  $X$  in the presence of  $Z$ , but not in its absence.

We now provide one real example from the literature. Most of the examples that we found are about indirect causality and use auxiliary variables that come from a pre-existing economic theory.

**Example 3** [*Indirect causality from money/credit to income*] Fackler (1985) found that neither money nor credit directly cause real output, but these variables play an indirect role in income determination. He also found that interest rates provide the link between the financial and real sectors, thus they can be viewed as the auxiliary variables that transmit the indirect causality from money/credit to income. In Section 8, we use the statistical procedures that we propose in this paper to re-examine these findings and check if interest rates are effectively the appropriate auxiliary variables.

The examples that we found in the literature show that an economic theory is needed in order to identify the appropriate auxiliary variables. In the absence of such theory, it seems difficult to find these variables that are needed to study the indirect/spurious causality. Hence the importance of our approach, which does not require knowledge of auxiliary variables.

## 4 Testing for indirect and spurious causalities

### 4.1 Testing for indirect causality

Definition 1 shows that there are three conditions that must be satisfied in order to have an indirect causality from  $Y$  to  $X$ . The first one [condition (i)] is simple to test as it only involves the observed variables  $X$  and  $Y$ . However, the other two conditions [conditions (ii) and (iii)] are difficult because the auxiliary variable  $Z$  is unknown, thus not observed. In the following, we propose to use factor analysis to identify  $Z$ . In particular, we use as a proxy of  $Z$  the factor(s) that we extract from a big data using principal component analysis.

Formally, condition (i) can be checked using the regression

$$X_{t+1} = \mu + \sum_{i=1}^p \beta_i X_{t+1-i} + \sum_{j=1}^q \alpha_j Y_{t+1-j} + \varepsilon_{t+1} \quad (9)$$

and an  $F$ -test for testing the null hypothesis

$$H_0 : \alpha_1 = \dots = \alpha_q = 0 \quad \text{vs} \quad H_1 : \text{No } H_0.$$

If  $H_0$  is rejected and  $Z$  is observed, we proceed to verify the condition (ii) using the regression

$$X_{t+1} = \eta + \sum_{i=1}^{\bar{p}} \gamma_i X_{t+1-i} + \sum_{j=1}^{\bar{q}} \lambda_j Y_{t+1-j} + \sum_{l=1}^{\bar{h}} \theta_l Z_{t+1-l} + e_{t+1} \quad (10)$$



and an  $F$ -test for testing the null hypothesis

$$\bar{H}_0 : \lambda_1 = \dots = \lambda_{\bar{q}} = 0 \quad \text{vs} \quad \bar{H}_1 : \text{No } \bar{H}_0.$$

If  $\bar{H}_0$  is not rejected, we proceed to check the condition (iii) using the regressions

$$Z_{t+1} = \nu + \sum_{i=1}^{\dot{p}} \kappa_i X_{t+1-i} + \sum_{j=1}^{\dot{q}} \psi_j Y_{t+1-j} + \sum_{l=1}^{\dot{h}} \rho_l Z_{t+1-l} + u_{t+1}, \quad (11)$$

$$X_{t+1} = \varpi + \sum_{i=1}^{\ddot{p}} \xi_i X_{t+1-i} + \sum_{j=1}^{\ddot{q}} \delta_j Y_{t+1-j} + \sum_{l=1}^{\ddot{h}} \varsigma_l Z_{t+1-l} + \epsilon_{t+1}, \quad (12)$$

and the  $F$ -tests for testing the null hypotheses

$$\begin{aligned} \dot{H}_0 : \psi_1 = \dots = \psi_{\dot{q}} = 0 \quad \text{vs} \quad \dot{H}_1 : \text{No } \dot{H}_0, \\ \ddot{H}_0 : \varsigma_1 = \dots = \varsigma_{\ddot{h}} = 0 \quad \text{vs} \quad \ddot{H}_1 : \text{No } \ddot{H}_0. \end{aligned}$$

If both  $\dot{H}_0$  and  $\ddot{H}_0$  are rejected, we conclude that  $Y$  indirectly causes  $X$ .

In practice, however,  $Z$  is not observed but it can be proxied by the factors that we extract from a big data that contains all economic variables that are available to econometricians. Formally, we consider an  $N$ -dimensional vector of large number of economic time series  $w_t = (w_{1t}, \dots, w_{Nt})^\top$  observed at each time  $t$ . We denote by  $W = (w_1, \dots, w_T)^\top$  the  $(T \times N)$ -dimensional matrix in which the  $t$ -th row is given by  $w_t$ . We assume that  $k$  common factors  $f_t$  are associated with the  $N$ -dimensional vector  $w_t$  according to the following equation:

$$w_t = \Lambda f_t + \varepsilon_t, \quad (13)$$

where  $\Lambda$  is an  $(N \times k)$ -dimensional matrix of factor loadings and  $\varepsilon_t$ 's are vectors of idiosyncratic shocks that could be cross-sectionally and temporally dependent.

To extract  $f$  from the  $(T \times N)$ -dimensional matrix  $W$ , we consider the factor model in Equation (13), which associates the  $N$ -dimensional vector  $w_t$  with the  $k$  common factors  $f_t$ . However, we remind the reader that the factors  $f_t$  and loadings  $\Lambda$  are not identified simultaneously since

$$w_t = \Lambda f_t + \varepsilon_t = \Lambda \Delta^{-1} \Delta f_t + \varepsilon_t = \Lambda^* f_t^* + \varepsilon_t, \quad (14)$$

for  $\Lambda^* = \Lambda \Delta^{-1}$ ,  $f_t^* = \Delta f_t$ , and  $\Delta$  a  $(k \times k)$ -dimensional positive definite matrix. Thus, we will only estimate the space spanned by the true factors and not the factors themselves. Simultaneous identification of the factors is, however, not essential for the statistical procedures that we propose to test for indirect/spurious causality. In other words, we only need to control for all available variables

and it doesn't matter how they are weighted in the factors. Using Equation (13) and each element  $w_{jt}$  of the vector  $w_t$ , the factor model is given by

$$w_{jt} = \vartheta_j^\top f_t + \varepsilon_{jt}, \text{ for } j = 1, \dots, N,$$

where  $\vartheta_j$  is a  $k$ -dimensional vector of factor loadings given by the  $j$ -th row of the matrix  $\Lambda$ . The factors  $f_t$  will be extracted using the principal component analysis (PCA) based on the following nonlinear least squares objective function

$$V(\tilde{f}, \tilde{\Lambda}) = \frac{1}{NT} \sum_{j=1}^N \sum_{t=1}^T \left( w_{j,t} - \tilde{\vartheta}_j^\top \tilde{f}_t \right)^2. \quad (15)$$

The function  $V(\tilde{f}, \tilde{\Lambda})$  depends on the hypothetical values of the factors  $\tilde{f} = (\tilde{f}_1, \dots, \tilde{f}_T)^\top$  and factor loadings  $\tilde{\Lambda} = (\tilde{\vartheta}_1, \dots, \tilde{\vartheta}_N)^\top$ . Let  $\hat{f}$  and  $\hat{\Lambda}$  be the minimizers of  $V(\tilde{f}, \tilde{\Lambda})$ . After concentrating out  $\tilde{f}$ , minimizing  $V(\tilde{f}, \tilde{\Lambda})$  is equivalent to maximizing  $\text{tr}(\tilde{\Lambda}^\top Y^\top Y \tilde{\Lambda})$  subject to  $\frac{\tilde{\Lambda}^\top \tilde{\Lambda}}{N} = I_k$ , where  $\text{tr}(\cdot)$  denotes the trace of a matrix and  $I_k$  is a  $(k \times k)$ -dimensional identity matrix. This represents the classical principal components problem that can be solved by setting the columns of  $\hat{\Lambda}$  to be equal to the eigenvectors of  $W^\top W$  corresponding to the  $k$  largest eigenvalues. The resulting principal components estimator of the matrix of the factors  $f = (f_1, \dots, f_T)^\top$  is:

$$\hat{f} = (\hat{f}_1, \dots, \hat{f}_T)^\top = \frac{W \hat{\Lambda}}{N}. \quad (16)$$

The computation of  $\hat{f}$  requires the calculation of the eigenvectors of the  $N \times N$  matrix  $W^\top W$  for  $N > T$ . Under some regularity conditions, Bai and Ng (2002) show that  $\hat{f}$  is a consistent estimator of  $f$ ; see Theorem 1 of Bai and Ng (2002).

We now replace  $Z$  in equations (10)-(12) by the extracted factor  $\hat{f}$ . Thus, to check condition (ii) of Definition 1, we use the feasible regression

$$X_{t+1} = \eta + \sum_{i=1}^{\bar{p}} \gamma_i X_{t+1-i} + \sum_{j=1}^{\bar{q}} \lambda_j Y_{t+1-j} + \sum_{l=1}^{\bar{h}} \theta_l \hat{f}_{t+1-l} + e_{t+1} \quad (17)$$

and the  $F$ -test for testing

$$\bar{H}_0 : \lambda_1 = \dots = \lambda_{\bar{q}} = 0 \quad \text{vs} \quad \bar{H}_1 : \text{No } \bar{H}_0. \quad (18)$$

Similarly, to verify condition (iii), we use the feasible regressions

$$\hat{f}_{t+1} = \nu + \sum_{i=1}^{\dot{p}} \kappa_i X_{t+1-i} + \sum_{j=1}^{\dot{q}} \psi_j Y_{t+1-j} + \sum_{l=1}^{\dot{h}} \rho_l \hat{f}_{t+1-l} + u_{t+1}, \quad (19)$$

$$X_{t+1} = \varpi + \sum_{i=1}^{\ddot{p}} \xi_i X_{t+1-i} + \sum_{j=1}^{\ddot{q}} \delta_j Y_{t+1-j} + \sum_{l=1}^{\ddot{h}} \varsigma_l \hat{f}_{t+1-l} + \epsilon_{t+1} \quad (20)$$

and the  $F$ -tests for testing the null hypotheses:

$$\dot{H}_0 : \psi_1 = \dots = \psi_{\dot{q}} = 0 \text{ vs } \dot{H}_1 : \text{No } \dot{H}_0 \quad (21)$$

$$\ddot{H}_0 : \varsigma_1 = \dots = \varsigma_{\ddot{h}} = 0 \text{ vs } \ddot{H}_1 : \text{No } \ddot{H}_0. \quad (22)$$

If both  $\dot{H}_0$  and  $\ddot{H}_0$  are rejected, we conclude that  $Y$  indirectly causes  $X$ .

## 4.2 Testing for spurious causality

To test for the spurious causality from  $Y$  to  $X$ , we need to check the conditions of Definition 2. For the spurious causality of type I, we have to check if conditions 1.(i) -1.(iii) are satisfied. To test condition 1.(i), we use the feasible regression

$$X_{t+1} = \mu + \sum_{i=1}^p \beta_i X_{t+1-i} + \sum_{j=1}^q \alpha_j Y_{t+1-j} + \sum_{l=1}^k \pi_j \hat{f}_{t+1-l} + \varepsilon_{t+1}, \quad (23)$$

and an  $F$ -test for testing the null hypothesis

$$H_0^{(1)} : \alpha_1 = \dots = \alpha_q = 0 \text{ vs } H_1^{(1)} : \text{No } H_0^{(1)}. \quad (24)$$

If  $H_0^{(1)}$  is rejected, we proceed to test condition 1.(ii) using the regression

$$X_{t+1} = \eta + \sum_{i=1}^{\bar{p}} \beta_i X_{t+1-i} + \sum_{j=1}^{\bar{q}} \alpha_j Y_{t+1-j} + e_{t+1},$$

and an  $F$ -test for testing the null hypothesis

$$\bar{H}_0^{(1)} : \alpha_1 = \dots = \alpha_{\bar{q}} = 0 \text{ vs } \bar{H}_1^{(1)} : \text{No } \bar{H}_0^{(1)}. \quad (25)$$

If  $\bar{H}_0^{(1)}$  is not rejected, we proceed to check the condition 1.(iii) before deciding about the presence of spurious causality of type I. Condition 1.(iii) can be verified using the feasible regressions:

$$\hat{f}_{t+1} = \nu + \sum_{i=1}^{\hat{p}} \kappa_i X_{t+1-i} + \sum_{j=1}^{\hat{q}} \psi_j Y_{t+1-j} + \sum_{l=1}^{\hat{h}} \rho_j \hat{f}_{t+1-l} + u_{t+1}, \quad (26)$$

$$X_{t+1} = \varpi + \sum_{i=1}^{\hat{p}} \xi_i X_{t+1-i} + \sum_{j=1}^{\hat{q}} \delta_j Y_{t+1-j} + \sum_{l=1}^{\hat{h}} \varsigma_j \hat{f}_{t+1-l} + \epsilon_{t+1}, \quad (27)$$

and the  $F$ -tests for testing the null hypotheses:

$$\dot{H}_0^{(1)} : \psi_1 = \dots = \psi_{\dot{q}} = 0 \text{ vs } \dot{H}_1^{(1)} : \text{No } \dot{H}_0^{(1)}, \quad (28)$$

$$\ddot{H}_0^{(1)} : \varsigma_1 = \dots = \varsigma_{\ddot{h}} = 0 \text{ vs } \ddot{H}_1^{(1)} : \text{No } \ddot{H}_0^{(1)}. \quad (29)$$

If both  $\dot{H}_0^{(1)}$  and  $\ddot{H}_0^{(1)}$  are rejected, we conclude that  $Y$  spuriously (type I) causes  $X$ .

For the spurious causality of type II, we need to check conditions 2.(i)-2.(iii). To test condition 2.(i), we use the regression

$$X_{t+1} = \mu + \sum_{i=1}^p \beta_i X_{t+1-i} + \sum_{j=1}^q \alpha_j Y_{t+1-j} + \varepsilon_{t+1},$$

and an  $F$ -test for testing the null hypothesis

$$H_0^{(2)} : \alpha_1 = \dots = \alpha_q = 0 \quad \text{vs} \quad H_1^{(2)} : \text{No } H_0^{(2)}.$$

If  $H_0^{(2)}$  is rejected, we proceed to test condition 2.(ii) using the feasible regression

$$X_{t+1} = \eta + \sum_{i=1}^{\bar{p}} \gamma_i X_{t+1-i} + \sum_{j=1}^{\bar{q}} \lambda_j Y_{t+1-j} + \sum_{l=1}^{\bar{h}} \theta_l \hat{f}_{t+1-l} + e_{t+1}$$

and an  $F$ -test for testing the null hypothesis

$$\bar{H}_0^{(2)} : \lambda_1 = \dots = \lambda_{\bar{q}} = 0 \quad \text{vs} \quad \bar{H}_1^{(2)} : \text{No } \bar{H}_0^{(2)}. \quad (30)$$

Thereafter, if  $\bar{H}_0^{(2)}$  is not rejected, we proceed to check the condition 2.(iii) before deciding about the presence of spurious causality of type II. Condition 2.(iii) can be verified using the feasible regressions

$$\begin{aligned} Y_{t+1} &= \nu + \sum_{i=1}^{\dot{p}} \kappa_i X_{t+1-i} + \sum_{j=1}^{\dot{q}} \psi_j Y_{t+1-j} + \sum_{l=1}^{\dot{h}} \rho_j \hat{f}_{t+1-l} + u_{t+1}, \\ X_{t+1} &= \varpi + \sum_{i=1}^{\ddot{p}} \xi_i X_{t+1-i} + \sum_{j=1}^{\ddot{q}} \delta_j Y_{t+1-j} + \sum_{l=1}^{\ddot{h}} \varsigma_j \hat{f}_{t+1-l} + \epsilon_{t+1}, \end{aligned}$$

and the  $F$ -tests for testing the null hypotheses:

$$\begin{aligned} \dot{H}_0^{(2)} : \rho_1 = \dots = \rho_{\dot{h}} = 0 \quad \text{vs} \quad \dot{H}_1^{(2)} : \text{No } \dot{H}_0^{(2)}, \\ \ddot{H}_0^{(2)} : \varsigma_1 = \dots = \varsigma_{\ddot{h}} = 0 \quad \text{vs} \quad \ddot{H}_1^{(2)} : \text{No } \ddot{H}_0^{(2)}. \end{aligned} \quad (31)$$

If both  $\dot{H}_0^{(2)}$  and  $\ddot{H}_0^{(2)}$  are rejected, we conclude that  $Y$  spuriously (type II) causes  $X$ .

We next derive the asymptotic distributions of the above non-causality tests in the presence of an estimated index  $\hat{f}$ , which we use as a proxy of the auxiliary variable  $Z$ .

## 5 Asymptotic distribution

In this section, we study the properties of indirect/spurious causality tests described in Section 4. In particular, we provide their asymptotic distributions in the presence of a consistent estimator of factors  $f$ . The assumptions required for the consistency of  $f$  are given by the following conditions.

**Assumption A:** For a positive generic constant  $\delta$ , we assume that: **(A1)**  $E \|f_t\|^4 \leq \delta < \infty$  and  $T^{-1} \sum_{t=1}^T f_t f_t^\top \xrightarrow{p} \Sigma_f > 0$ , where  $\Sigma_f$  is a  $(k \times k)$  non-random positive definite matrix and  $\|\cdot\|$  denotes the Euclidean norm; **(A2)** If  $\Lambda$  is deterministic, then  $\|\lambda_j\| \leq \delta < \infty$ , where  $\lambda_j$  is the  $j$ -th row of the factor loadings matrix  $\Lambda$ . If it is stochastic, then  $E \|\lambda_j\|^4 \leq \delta$ . Furthermore,  $N^{-1} \Lambda^\top \Lambda \xrightarrow{p} \Sigma_\Lambda > 0$ , as  $N \rightarrow \infty$ , where  $\Sigma_\Lambda$  is a  $(k \times k)$  non-random matrix; **(A3)** For all  $N$  and  $T$ , we have: **(i)**  $E(\varepsilon_{it}) = 0$  and  $E|\varepsilon_{it}|^8 \leq \delta$ ; **(ii)** For  $N^{-1} E(\varepsilon_s^\top \varepsilon_t) = \gamma_N(s, t)$ , we have  $|\gamma_N(s, s)| \leq \delta$ , where  $s = 1, \dots, T$ , and  $T^{-1} \sum_{1 \leq s, t \leq N} |\gamma_N(s, t)| \leq \delta$ ; **(iii)** For  $E(\varepsilon_{it} \varepsilon_{jt}) = \tau_{ij, t}$ , we have  $|\tau_{ij, t}| \leq \tau_{ij}$ ,  $\forall t$ , and  $N^{-1} \sum_{1 \leq i, j \leq N} |\tau_{ij}| \leq \delta$ ; **(iv)** For  $E(\varepsilon_{it} \varepsilon_{js}) = \tau_{ij, ts}$ , we have  $(TN)^{-1} \sum_{1 \leq i, j, s, t \leq N} |\tau_{ij, ts}| \leq \delta$ ; and **(v)**  $E \left| N^{-\frac{1}{2}} \sum_{i=1}^N [\varepsilon_{it} \varepsilon_{js} - E(\varepsilon_{it} \varepsilon_{js})] \right|^4 \leq \delta$ ,  $\forall s, t$ ; and **(A4)**  $E \left( N^{-1} \sum_{i=1}^N \left\| T^{-1} \sum_{t=1}^T f_t \varepsilon_{it} \right\|^2 \right) \leq \delta$ .

Assumptions **(A1)** and **(A2)** are standard in the literature on factor analysis; see Stock and Watson (2002), Bai and Ng (2002, 2006, 2008) and Bai (2003). They represent moment conditions on factors  $f_t$  and factor loadings  $\Lambda$ , and they ensure that the factors are non-degenerate and their contribution to the variance of the data is nontrivial. Assumptions **(A3)-(i)** to **(A3)-(v)** allow for heteroskedasticity and weak correlation between the components of the vector of idiosyncratic shocks  $\varepsilon_t$  in (13). Under these assumptions both cross-sectional and serial correlations are allowed. We next establish the asymptotic distributions of indirect/spurious causality tests in Section 4.

## 5.1 Indirect causality

We derive the asymptotic distributions of tests of conditions of indirect causality in Definition 1. We focus on testing conditions (ii) and (iii), because the test of condition (i) depends only on the observed variables. We introduce the following notations. Define the vector of parameters  $\tau_1^{Ind} = (\eta, \gamma', \lambda', \theta)'$ , where  $\gamma = (\gamma_1, \dots, \gamma_{\bar{p}})'$ , and  $\lambda$  and  $\theta$  can be defined in similar way. Let  $\hat{z}_{1t} = (1, X_t, \dots, X_{t+1-\bar{p}}, Y_t, \dots, Y_{t+1-\bar{q}}, \hat{f}_t, \dots, \hat{f}_{t+1-\bar{h}})'$  and  $\hat{\tau}_1^{Ind} = (\hat{\eta}, \hat{\gamma}', \hat{\lambda}', \hat{\theta})'$ , where  $\hat{\tau}_1^{Ind}$  is the least squares estimates obtained from the regression of  $X_{t+1}$  on the constant,  $(X_t, \dots, X_{t+1-\bar{p}})'$ ,  $(Y_t, \dots, Y_{t+1-\bar{q}})'$ , and the estimated factors  $(\hat{f}_t, \dots, \hat{f}_{t+1-\bar{h}})'$ . Henceforth,  $0_{m,n}$  is the  $m \times n$  matrix of zeros and  $I_n$  is the  $n \times n$  identity matrix.

First, we test condition (ii) using the following  $F$ -test statistic, which tests that all the coefficients of  $Y_t, \dots, Y_{t+1-\bar{q}}$  in the regression (17) are jointly equal to zero:

$$F_T^{Ind, \lambda} = \left( \sqrt{T} R^{Ind, \lambda} \hat{\tau}_1^{Ind} \right) \left( R^{Ind, \lambda} \hat{\Sigma}_{\tau_1^{Ind}} R^{Ind, \lambda} \right)^{-1} \left( \sqrt{T} R^{Ind, \lambda} \hat{\tau}_1^{Ind} \right)', \quad (32)$$

where the selection matrix  $R^{Ind,\lambda} = (0_{\bar{q},1+\bar{p}}, I_{\bar{q}}, 0_{\bar{q},\bar{h}})$  and

$$\hat{\Sigma}_{\tau_1^{Ind}} = \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{1t} \hat{z}'_{1t} \right)^{-1} \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{e}_{t+1}^2 \hat{z}_{1t} \hat{z}'_{1t} \right) \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{1t} \hat{z}'_{1t} \right)^{-1}, \quad (33)$$

with  $\hat{e}_{t+1} = X_{t+1} - \hat{z}'_{1t} \hat{\tau}_1^{Ind}$  are the residuals. The following theorem demonstrates that the  $F_T^{Ind,\lambda}$  test statistic is asymptotically distributed as a chi-squared distribution with  $\bar{q}$  degrees of freedom [see the proof of Theorem 1 in Appendix A].

**Theorem 1 :** *Let Assumption A hold. Under the null hypothesis (18), if  $\sqrt{T}/N \rightarrow 0$ , then we have*

$$F_T^{Ind,\lambda} \rightarrow_d \chi_{\bar{q}}^2,$$

where the  $F_T^{Ind,\lambda}$  statistic is defined in (32).

Theorem 1 is stated and proved under the general case of heteroskedasticity. However, the above result is still valid under homoskedasticity, with a consistent estimator of variance-covariance matrix:

$$\hat{\Sigma}_{\tau_1^{Ind}}^* = \hat{\sigma}_e^2 \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{1t} \hat{z}'_{1t} \right)^{-1}, \quad (34)$$

where  $\hat{\sigma}_e^2 = T^{-1} \sum_{t=1}^{T-1} \hat{e}_{t+1}^2$ . It can be shown that the difference between the two estimators in (33) and (34) is asymptotically negligible under homoskedasticity. Furthermore, the proof of Theorem 1 indicates that for  $\sqrt{T}/N \rightarrow 0$ , having estimated factors as regressors does not affect the root- $T$  consistency of the least squares estimates of  $\tau_1^{Ind}$ , except that the variance-covariance matrix  $\Sigma_{\tau_1^{Ind}}$  will be different and can be consistently estimated by  $\hat{\Sigma}_{\tau_1^{Ind}}$ . However, if  $\sqrt{T}/N \rightarrow c > 0$ , then there are two additional terms that do not vanish asymptotically, thus  $\hat{\tau}_1^{Ind}$  will have an asymptotic bias term. For more details the reader is referred to the proof of Lemma 3 in Appendix A.

We now provide the test statistics that one can use to test condition (iii) or equivalently the null hypotheses (21) and (22). The latter hypotheses can be tested using the following  $F_T^{Ind,\psi}$  and  $F_T^{Ind,s}$  test statistics that test if all the coefficients of the vectors  $(Y_t, \dots, Y_{t+1-\hat{q}})'$  and  $(\hat{f}_t, \dots, \hat{f}_{t+1-\hat{h}})'$  in the regressions (19) and (20) are jointly equal to zero, respectively. The  $F_T^{Ind,\psi}$  test statistic for testing the null hypothesis (21) is given by:

$$F_T^{Ind,\psi} = \left( \sqrt{T} R^{Ind,\psi} \hat{\tau}_2^{Ind} \right) \left( R^{Ind,\psi} \hat{\Sigma}_{\tau_2^{Ind}} R^{Ind,\psi'} \right)^{-1} \left( \sqrt{T} R^{Ind,\psi} \hat{\tau}_2^{Ind} \right)', \quad (35)$$

where  $\hat{\tau}_2^{Ind} = (\hat{\nu}, \hat{\kappa}', \hat{\psi}', \hat{\rho}')'$ , the selection matrix  $R^{Ind,\psi} = (0_{\hat{q},1+\hat{p}}, I_{\hat{q}}, 0_{\hat{q},\hat{h}})$ , and

$$\hat{\Sigma}_{\tau_2^{Ind}} = \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{2t} \hat{z}'_{2t} \right)^{-1} \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{u}_{t+1}^2 \hat{z}_{2t} \hat{z}'_{2t} \right) \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{2t} \hat{z}'_{2t} \right)^{-1},$$

with  $\hat{u}_{t+1} = \hat{f}_{t+1} - \hat{z}'_{2t} \hat{\tau}_2^{Ind}$  are the residuals from the regression in (19) and

$$\hat{z}_{2t} = (1, X_t, \dots, X_{t+1-p}, Y_t, \dots, Y_{t+1-q}, \hat{f}_t, \dots, \hat{f}_{t+1-h})'.$$

Similarly, the  $F_T^{Ind, \varsigma}$  test statistic for testing the null hypothesis (22) is given by:

$$F_T^{Ind, \varsigma} = \left( \sqrt{T} R^{Ind, \varsigma} \hat{\tau}_3^{Ind} \right) \left( R^{Ind, \varsigma} \hat{\Sigma}_{\tau_3^{Ind}} R^{Ind, \varsigma'} \right)^{-1} \left( \sqrt{T} R^{Ind, \varsigma} \hat{\tau}_3^{Ind} \right)', \quad (36)$$

where  $\hat{\tau}_3^{Ind} = (\hat{\omega}, \hat{\xi}', \hat{\delta}', \hat{\zeta}')'$ , the selection matrix  $R^{Ind, \varsigma} = (0_{\check{h}, 1+\check{p}+\check{q}}, I_{\check{h}})$ , and

$$\hat{\Sigma}_{\tau_3^{Ind}} = \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{3t} \hat{z}'_{3t} \right)^{-1} \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{\epsilon}_{t+1}^2 \hat{z}_{3t} \hat{z}'_{3t} \right) \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{3t} \hat{z}'_{3t} \right)^{-1},$$

with  $\hat{\epsilon}_{t+1} = X_{t+1} - \hat{z}'_{3t} \hat{\tau}_3^{Ind}$  are the residuals from the regression in (20) and

$$\hat{z}_{3t} = (1, X_t, \dots, X_{t+1-p}, Y_t, \dots, Y_{t+1-q}, \hat{f}_t, \dots, \hat{f}_{t+1-h})'.$$

The following theorem shows that the test statistics  $F_T^{Ind, \psi}$  and  $F_T^{Ind, \varsigma}$  are asymptotically distributed as chi-squared distributions with  $\check{q}$  and  $\check{h}$  degrees of freedom, respectively [see the proof of Theorem 2 in Appendix A].

**Theorem 2 :** *Let Assumption A hold. Under the null hypotheses (21) and (22), if  $\sqrt{T}/N \rightarrow 0$ , then we have*

$$F_T^{Ind, \psi} \rightarrow_d \chi_{\check{q}}^2 \text{ and } F_T^{Ind, \varsigma} \rightarrow_d \chi_{\check{h}}^2,$$

where the  $F_T^{Ind, \psi}$  and  $F_T^{Ind, \varsigma}$  statistics are defined in (35) and (36), respectively.

Similar to Theorem 1, the result in Theorem 2 is still valid under homoskedasticity, with consistent estimators of the variance-covariance matrices  $\Sigma_{\tau_2^{Ind}}$  and  $\Sigma_{\tau_3^{Ind}}$ :

$$\hat{\Sigma}_{\tau_2^{Ind}}^* = \hat{\sigma}_u^2 \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{2t} \hat{z}'_{2t} \right)^{-1} \text{ and } \hat{\Sigma}_{\tau_3^{Ind}}^* = \hat{\sigma}_\epsilon^2 \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{3t} \hat{z}'_{3t} \right)^{-1},$$

respectively, where  $\hat{\sigma}_u^2 = (1/T) \sum_{t=1}^{T-1} \hat{u}_{t+1}^2$  and  $\hat{\sigma}_\epsilon^2 = (1/T) \sum_{t=1}^{T-1} \hat{\epsilon}_{t+1}^2$ . Furthermore, we can show that the difference between the estimators  $\hat{\Sigma}_{\tau_2^{Ind}}$  and  $\hat{\Sigma}_{\tau_2^{Ind}}^*$  (resp.  $\hat{\Sigma}_{\tau_3^{Ind}}$  and  $\hat{\Sigma}_{\tau_3^{Ind}}^*$ ) is asymptotically negligible under homoskedasticity. Moreover, our results show that for  $\sqrt{T}/N \rightarrow 0$ , having estimated factors as regressand or regressors does not affect the root- $T$  consistency of the least squares estimates of  $\tau_2^{Ind}$  and  $\tau_3^{Ind}$ , except that the variance-covariance matrices  $\Sigma_{\tau_2^{Ind}}$  and  $\Sigma_{\tau_3^{Ind}}$  will have different expressions and can be consistently estimated by  $\hat{\Sigma}_{\tau_2^{Ind}}$  and  $\hat{\Sigma}_{\tau_3^{Ind}}$ , respectively. However, if  $\sqrt{T}/N \rightarrow c > 0$ , then there are additional terms that do not vanish asymptotically, thus  $\hat{\tau}_2^{Ind}$  and  $\hat{\tau}_3^{Ind}$  will have asymptotic bias terms.

## 5.2 Spurious causality

We now study the properties of tests of conditions of spurious causality in Definition 2. For type I spurious causality, we focus on providing the asymptotic distributions of tests of conditions (i) and (iii), since the test of condition (ii) only depends on observed variables and can be tested using the standard  $F$ -test. For type II spurious causality, we only derive the asymptotic distributions of tests of conditions (ii) and (iii), again because the test of condition (i) involves observed variables only.

Regarding the type I spurious causality, conditions (i), (iii)-(a) and (iii)-(b) can be tested using the following  $F$ -test statistics

$$\begin{aligned} F_T^{SI,\alpha} &= \left( \sqrt{T} R^{SI,\alpha} \hat{\tau}_1^{si} \right) \left( R^{SI,\alpha} \hat{\Sigma}_{\hat{\tau}_1^{si}} R^{SI,\alpha'} \right)^{-1} \left( \sqrt{T} R^{SI,\alpha} \hat{\tau}_1^{si} \right)', \\ F_T^{SI,\psi} &= \left( \sqrt{T} R^{SI,\psi} \hat{\tau}_2^{si} \right) \left( R^{SI,\psi} \hat{\Sigma}_{\hat{\tau}_2^{si}} R^{SI,\psi'} \right)^{-1} \left( \sqrt{T} R^{SI,\psi} \hat{\tau}_2^{si} \right)', \\ F_T^{SI,\varsigma} &= \left( \sqrt{T} R^{SI,\varsigma} \hat{\tau}_3^{si} \right) \left( R^{SI,\varsigma} \hat{\Sigma}_{\hat{\tau}_3^{si}} R^{SI,\varsigma'} \right)^{-1} \left( \sqrt{T} R^{SI,\varsigma} \hat{\tau}_3^{si} \right)', \end{aligned} \quad (37)$$

respectively, where  $\hat{\tau}_1^{si} = (\hat{\mu}, \hat{\beta}', \hat{\alpha}', \hat{\pi}')'$ ,  $\hat{\tau}_2^{si} = (\hat{\nu}, \hat{\kappa}', \hat{\psi}', \hat{\rho}')'$ ,  $\hat{\tau}_3^{si} = (\hat{\omega}, \hat{\xi}', \hat{\delta}', \hat{\zeta}')'$ , the selection matrices  $R^{SI,\alpha} = (0_{q,1+p}, I_q, 0_{q,k})$ ,  $R^{SI,\psi} = (0_{\dot{q},1+\dot{p}}, I_{\dot{q}}, 0_{\dot{q},\dot{h}})$ ,  $R^{SI,\varsigma} = (0_{\ddot{h},1+\ddot{p}+\ddot{q}}, I_{\ddot{h}})$ , and

$$\begin{aligned} \hat{\Sigma}_{\hat{\tau}_1^{si}} &= \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{1t}^{si} \hat{z}_{1t}^{si'} \right)^{-1} \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{\epsilon}_{t+1}^2 \hat{z}_{1t}^{si} \hat{z}_{1t}^{si'} \right) \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{1t}^{si} \hat{z}_{1t}^{si'} \right)^{-1} \\ \hat{\Sigma}_{\hat{\tau}_2^{si}} &= \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{2t}^{si} \hat{z}_{2t}^{si'} \right)^{-1} \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{u}_{t+1}^2 \hat{z}_{2t}^{si} \hat{z}_{2t}^{si'} \right) \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{2t}^{si} \hat{z}_{2t}^{si'} \right)^{-1} \\ \hat{\Sigma}_{\hat{\tau}_3^{si}} &= \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{3t}^{si} \hat{z}_{3t}^{si'} \right)^{-1} \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{\epsilon}_{t+1}^2 \hat{z}_{3t}^{si} \hat{z}_{3t}^{si'} \right) \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{3t}^{si} \hat{z}_{3t}^{si'} \right)^{-1} \end{aligned}$$

with  $\hat{\epsilon}_{t+1} = X_{t+1} - \hat{z}_{1t}^{si'} \hat{\tau}_1^{si}$ ,  $\hat{u}_{t+1} = \hat{f}_{t+1} - \hat{z}_{2t}^{si'} \hat{\tau}_2^{si}$ ,  $\hat{\epsilon}_{t+1} = X_{t+1} - \hat{z}_{3t}^{si'} \hat{\tau}_3^{si}$ , and

$$\begin{aligned} \hat{z}_{1t}^{si} &= (1, X_t, \dots, X_{t+1-p}, Y_t, \dots, Y_{t+1-q}, \hat{f}_t, \dots, \hat{f}_{t+1-k})', \\ \hat{z}_{2t}^{si} &= (1, X_t, \dots, X_{t+1-\dot{p}}, Y_t, \dots, Y_{t+1-\dot{q}}, \hat{f}_t, \dots, \hat{f}_{t+1-\dot{h}})', \\ \hat{z}_{3t}^{si} &= (1, X_t, \dots, X_{t+1-\ddot{p}}, Y_t, \dots, Y_{t+1-\ddot{q}}, \hat{f}_t, \dots, \hat{f}_{t+1-\ddot{h}})'. \end{aligned}$$

The following theorem shows that the test statistics  $F_T^{SI,\lambda}$ ,  $F_T^{SI,\rho}$ , and  $F_T^{SI,\varsigma}$  are asymptotically distributed as chi-squared distributions with  $q$ ,  $\dot{q}$ , and  $\ddot{h}$  degrees of freedom, respectively. The proof of Theorem 3 is omitted since it is similar to those of Theorems 1 and 2.

**Theorem 3** : *Let Assumption A holds. Under the null hypotheses (24), (28), and (29), if  $\sqrt{T}/N \rightarrow 0$ , then we have*

$$F_T^{SI,\alpha} \rightarrow_d \chi_q^2, \quad F_T^{SI,\psi} \rightarrow_d \chi_{\dot{q}}^2, \quad \text{and} \quad F_T^{SI,\varsigma} \rightarrow_d \chi_{\ddot{h}}^2,$$

where the test statistics  $F_T^{SI,\alpha}$ ,  $F_T^{SI,\psi}$  and  $F_T^{SI,\varsigma}$  are defined in (37), respectively.



Concerning the type II spurious causality, conditions (ii), (iii)-(a) and (iii)-(b) can be tested using the following  $F$ -test statistics

$$\begin{aligned} F_T^{SII,\lambda} &= \left( \sqrt{T} R^{sii,\lambda} \hat{\tau}_1^{sii} \right) \left( R^{sii,\lambda} \hat{\Sigma}_{\hat{\tau}_1^{sii}} R^{sii,\lambda'} \right)^{-1} \left( \sqrt{T} R^{sii,\lambda} \hat{\tau}_1^{sii} \right)', \\ F_T^{SII,\rho} &= \left( \sqrt{T} R^{sii,\rho} \hat{\tau}_2^{sii} \right) \left( R^{sii,\rho} \hat{\Sigma}_{\hat{\tau}_2^{sii}} R^{sii,\rho'} \right)^{-1} \left( \sqrt{T} R^{sii,\rho} \hat{\tau}_2^{sii} \right)', \\ F_T^{SII,\varsigma} &= \left( \sqrt{T} R^{sii,\varsigma} \hat{\tau}_3^{sii} \right) \left( R^{sii,\varsigma} \hat{\Sigma}_{\hat{\tau}_3^{sii}} R^{sii,\varsigma'} \right)^{-1} \left( \sqrt{T} R^{sii,\varsigma} \hat{\tau}_3^{sii} \right)', \end{aligned} \quad (38)$$

respectively, where  $\hat{\tau}_1^{sii} = (\hat{\eta}, \hat{\gamma}', \hat{\lambda}', \hat{\theta}')'$ ,  $\hat{\tau}_2^{sii} = (\hat{\nu}, \hat{\kappa}', \hat{\psi}', \hat{\rho}')'$ ,  $\hat{\tau}_3^{sii} = (\hat{\omega}, \hat{\xi}', \hat{\delta}', \hat{\zeta}')'$ , the selection matrices  $R^{sii,\lambda} = (0_{\bar{q},1+\bar{p}}, I_{\bar{q}}, 0_{\bar{q},\bar{h}})$ ,  $R^{sii,\rho} = (0_{\bar{h},1+\bar{p}+\bar{q}}, I_{\bar{h}})$ ,  $R^{sii,\varsigma} = (0_{\bar{h},1+\bar{p}+\bar{q}}, I_{\bar{h}})$ , and

$$\begin{aligned} \hat{\Sigma}_{\hat{\tau}_1^{sii}} &= \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{1t}^{sii} \hat{z}_{1t}^{siii'} \right)^{-1} \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{e}_{t+1}^2 \hat{z}_{1t}^{sii} \hat{z}_{1t}^{siii'} \right) \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{1t}^{sii} \hat{z}_{1t}^{siii'} \right)^{-1} \\ \hat{\Sigma}_{\hat{\tau}_2^{sii}} &= \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{2t}^{sii} \hat{z}_{2t}^{siii'} \right)^{-1} \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{u}_{t+1}^2 \hat{z}_{2t}^{sii} \hat{z}_{2t}^{siii'} \right) \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{2t}^{sii} \hat{z}_{2t}^{siii'} \right)^{-1} \\ \hat{\Sigma}_{\hat{\tau}_3^{sii}} &= \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{3t}^{sii} \hat{z}_{3t}^{siii'} \right)^{-1} \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{e}_{t+1}^2 \hat{z}_{3t}^{sii} \hat{z}_{3t}^{siii'} \right) \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_{3t}^{sii} \hat{z}_{3t}^{siii'} \right)^{-1} \end{aligned}$$

with  $\hat{e}_{t+1} = X_{t+1} - \hat{z}_{1t}^{siii'} \hat{\tau}_1^{sii}$ ,  $\hat{u}_{t+1} = Y_{t+1} - \hat{z}_{2t}^{siii'} \hat{\tau}_2^{sii}$ ,  $\hat{e}_{t+1} = X_{t+1} - \hat{z}_{3t}^{siii'} \hat{\tau}_3^{sii}$ , and

$$\begin{aligned} \hat{z}_{1t}^{sii} &= (1, X_t, \dots, X_{t+1-\bar{p}}, Y_t, \dots, Y_{t+1-\bar{q}}, \hat{f}_t, \dots, \hat{f}_{t+1-\bar{h}})', \\ \hat{z}_{2t}^{sii} &= (1, X_t, \dots, X_{t+1-\bar{p}}, Y_t, \dots, Y_{t+1-\bar{q}}, \hat{f}_t, \dots, \hat{f}_{t+1-\bar{h}})', \\ \hat{z}_{3t}^{sii} &= (1, X_t, \dots, X_{t+1-\bar{p}}, Y_t, \dots, Y_{t+1-\bar{q}}, \hat{f}_t, \dots, \hat{f}_{t+1-\bar{h}})'. \end{aligned}$$

The following theorem demonstrates that the test statistics  $F_T^{SII,\lambda}$ ,  $F_T^{SII,\rho}$ , and  $F_T^{SII,\varsigma}$  are asymptotically distributed as chi-squared distributions with  $\bar{q}$ ,  $\bar{h}$ , and  $\bar{h}$  degrees of freedom, respectively. The proof of Theorem 4 follows naturally in a similar way as the one of Theorem 1, and therefore it is omitted.

**Theorem 4** : *Let Assumption A hold. Under the null hypotheses (30) and (31), if  $\sqrt{T}/N \rightarrow 0$ , then we have*

$$F_T^{SII,\lambda} \rightarrow_d \chi_{\bar{q}}^2, \quad F_T^{SII,\rho} \rightarrow_d \chi_{\bar{h}}^2, \quad \text{and} \quad F_T^{SII,\varsigma} \rightarrow_d \chi_{\bar{h}}^2,$$

where the test statistics  $F_T^{SII,\lambda}$ ,  $F_T^{SII,\rho}$  and  $F_T^{SII,\varsigma}$  are defined in (38), respectively.

As in Section 5.1, the results in Theorems 3 and 4 show that for  $\sqrt{T}/N \rightarrow 0$ , having estimated factors as regressand or regressors does not affect the root- $T$  consistency of the least squares estimates of the coefficients used to test type I and type II spurious causalities. These results work under the general case of heteroskedasticity, and they are still valid under homoskedasticity.

## 6 Identification of the auxiliary variables

Unfortunately, the tests developed in the previous sections only detect the presence of indirect/spurious causality and cannot provide information about the nature of the auxiliary variables that transmit/induce these effects. In this section, we suggest a simple statistical procedure to identify the auxiliary variables that are responsible for the transmission/induction of indirect/spurious causality.

The literature on the interpretation of factors extracted using factor analysis suggests to use marginal regressions where each factor is regressed on each of the variables of the big data, see Ludvigson and Ng (2009) and the references therein. Thereafter, it uses the coefficient of determination  $R^2$  to order the variables according to their importance in terms of explaining each factor. Thus, the variable(s) that produce(s) high  $R^2_{(s)}$  are used to interpret the factor. Following this literature, we propose to identify the auxiliary variables in the following way:

1. Test the conditions of indirect/spurious causality using as an auxiliary variable  $f_t = w_{j,t}$ , for  $j = 1, \dots, N$ , where  $w_{j,t}$  is one of the variables of the big data  $W$  defined in Section 4.1;
2. Eliminate all  $w_{j,t}$ , for  $j = 1, \dots, N$ , which do not satisfy the conditions of indirect/spurious causality. We denote the subset of  $W$  with all variables satisfying the conditions of indirect/spurious causality by  $W^{sub} = \left\{ \left( w_t^{(1)}, \dots, w_t^{(\bar{N})} \right), \text{ for } t = 1, \dots, T \text{ and where } \bar{N} \leq N \right\}$ .  $W^{sub}$  can be viewed as a subset of auxiliary variables that transmit/induce indirect/spurious causality. If this subset is sufficiently large or because of possible interaction between the auxiliary variables, we can consider the next additional step;
3. Use the factor analysis where  $\bar{k}$  common factors  $f_t^{sub}$  are associated with the  $\bar{N}$ -dimensional vector  $w_t^{sub}$  according to the following equation:

$$w_t^{sub} = \Lambda^{sub} f_t^{sub} + \varepsilon_t,$$

where  $w_t^{sub} = (w_t^{(1)}, \dots, w_t^{(\bar{N})})^\top$ , for  $t = 1, \dots, T$ ,  $\Lambda^{sub}$  is an  $(\bar{N} \times \bar{k})$ -dimensional matrix of factor loadings and  $\varepsilon_t$ 's are vectors of idiosyncratic shocks that could be cross-sectionally and temporally dependent. We can then use the following marginal regressions, as in Ludvigson and Ng (2009), to interpret  $f_t^{sub}$ . For example, for one factor ( $f_t^{sub}$ ) model, we run the marginal regressions

$$f_t^{sub} = \alpha^{sub} + \beta^{sub} w_t^{(j)} + u_t$$

for each  $j = 1, \dots, \bar{N}$  and obtain the corresponding  $R^2$ s. Thus, the auxiliary factor  $f_t^{sub}$  can be interpreted in terms of the variables in  $W^{sub}$  that generate high  $R^2$ s.

## 7 Monte Carlo simulations

We conduct a Monte Carlo simulation to investigate the performance of statistical procedures proposed in the previous sections. Our primary focus is on assessing the empirical size and power of the tests in Theorems 1-3 under a variety of data-generating processes, different sample sizes and different numbers of auxiliary variables.

### 7.1 Indirect causality

We first describe the data-generating processes (DGPs) that we use in the simulations to assess the performance of the tests in Theorems 1 and 2. Initially, we consider a set of DGPs in which indirect causality is transmitted through one auxiliary variable  $Z_1$ , among a total of  $N$  variables  $\{Z_1, \dots, Z_N\}$  that represent the economy. In particular, we consider the following processes:

$$X_t = \mu_X^{(1)} + \phi_X^{(1)} X_{t-1} + \phi_Y^{(1)} Y_{t-1} + \varepsilon_{X,t}^{(1)}, \text{ with } \varepsilon_{X,t}^{(1)} \sim N(0, 1), \quad (39)$$

$$X_t = \mu_X^{(2)} + \phi_X^{(2)} X_{t-1} + \phi_Y^{(2)} Y_{t-2} + \phi_Z^{(2)} Z_{1,t-1} + \varepsilon_{X,t}^{(2)}, \text{ with } \varepsilon_{X,t}^{(2)} \sim N(0, 1), \quad (40)$$

$$Z_{1,t} = \mu_Z - \phi_Y^{(2)}/\phi_Z^{(2)} Y_{t-1} + \phi_Z Z_{1,t-1} + \varepsilon_{Z_1,t}, \text{ with } \varepsilon_{Z_1,t} \sim N(0, 1) \text{ and } \phi_Z^{(2)} \neq 0, \quad (41)$$

$$Z_{i,t} = \varepsilon_{Z_i,t} \sim N(0, 1), \text{ for } i = 2, \dots, N, \text{ with } \varepsilon_{Z_i,t} \text{ mutually independent,}$$

where the error terms  $\varepsilon_{X,t}^{(1)}$ ,  $\varepsilon_{X,t}^{(2)}$ , and  $\varepsilon_{Z_i,t}$  for  $i = 1, \dots, N$ , are assumed to be mutually independent and  $X$  has to satisfy both equations (39) and (40). Numerical values for the coefficients  $\mu_X^{(1)}$ ,  $\mu_X^{(2)}$ ,  $\mu_Z$ ,  $\phi_X^{(1)}$ ,  $\phi_X^{(2)}$ ,  $\phi_Y^{(1)}$ ,  $\phi_Y^{(2)}$ ,  $\phi_Z^{(2)}$ ,  $\phi_Z$  will be specified later. The functional forms of DGPs are selected such that there is an indirect causality from  $Y$  to  $X$ . According to Definition (1), indirect causality occurs if: (i)  $Y$  Granger causes  $X$  with respect to the information set  $I_X(t)$ ; (ii)  $Y$  does not Granger cause  $X$  with respect to the information set  $I(t) - I_Y(t)$ ; and (iii)  $Y$  Granger causes  $Z$  and  $Z$  Granger causes  $X$  with respect to the information sets  $I(t) - I_Y(t)$  and  $I(t) - I_Z(t)$ , respectively. Thus, condition (i) will be satisfied if we choose the coefficient  $\phi_Y^{(1)}$  to be different from zero. Furthermore, if we assume that  $\phi_Z^{(2)} \neq 0$  and  $\phi_Y^{(2)} \neq 0$ , then  $Z$  Granger causes  $X$  and  $Y$  Granger causes  $Z$ , respectively, hence condition (iii) is satisfied. What about condition (ii)? Combining equations (40) and (41) leads to

$$X_t = \left( \mu_X^{(2)} + \phi_Z^{(2)} \mu_Z \right) + \phi_X^{(2)} X_{t-1} + \phi_Z^{(2)} \phi_Z Z_{1,t-2} + \phi_Z^{(2)} \varepsilon_{Z_1,t-1} + \varepsilon_{X,t}^{(1)},$$

which indicates that  $Y$  does not cause  $X$  in the presence of  $Z$ , hence condition (ii) is also satisfied.

Depending on whether or not  $\phi_Y^{(1)}$ ,  $\phi_Z^{(2)}$  and  $\phi_Y^{(2)}$  are taken to be equal to zero, the following steps can be performed to simulate a sample of  $T$  observations on  $X$ ,  $Y$  and  $Z$  under the absence/presence of indirect causality from  $Y$  to  $X$  transmitted through the auxiliary variable  $Z_1$ :

Table 1: Data-generating processes: Direct causality

DGPs	Variables of interest		
	$X_t =$	$Y_t =$	$Z_{jt}$ for $j = 1, \dots, N$
DGP1	$0.5 + 0.2X_{t-1} + 0.3Y_{t-1} + \varepsilon_{1t}$	$0.5 + 0.5Y_{t-1} + \varepsilon_{2t}$	$Z_{1t} = \varepsilon_{3t}, Z_{jt} = \varepsilon_{j+2,t}$
DGP2	$0.5 + 0.2X_{t-1} + 0.7Y_{t-1} + \varepsilon_{1t}$	$0.5 + 0.5Y_{t-1} + \varepsilon_{2t}$	$Z_{1t} = \varepsilon_{3t}, Z_{jt} = \varepsilon_{j+2,t}$
DGP3	$0.5 + 0.2X_{t-1} + 0.3Y_{t-1} + \varepsilon_{1t}$	$0.5 + 0.5Y_{t-1} + \varepsilon_{2t}$	$Z_{1t} = 0.5 + 0.5Z_{1,t-1} + \varepsilon_{3t}, Z_{jt} = \varepsilon_{j+2,t}$
DGP4	$0.5 + 0.2X_{t-1} + 0.7Y_{t-1} + \varepsilon_{1t}$	$0.5 + 0.5Y_{t-1} + \varepsilon_{2t}$	$Z_{1t} = 0.5 + 0.5Z_{1,t-1} + \varepsilon_{3t}, Z_{jt} = \varepsilon_{j+2,t}$

**Note:** This table summarizes the DGPs, with a direct causality from  $Y$  to  $X$ , that we consider in the simulation study to examine the performance of the tests in Theorem 1 for testing the nulls outlined in (18), (21) and (22). The error terms  $\varepsilon_{it}$ , for  $i = 1, 2, 3$ , and  $\varepsilon_{j+2,t}$  for  $j = 2, \dots, N$  are  $N + 2$  mutually independent standard normal random variables, where  $N$  can be large indicating the richness of the data environment.

(1) Choose the initial values  $X_2$ ,  $Z_2$ , and  $Y_1$ , and generate  $X_3$  using Equation (40)

$$X_3 = \mu_X^{(2)} + \phi_X^{(2)}X_2 + \phi_Y^{(2)}Y_1 + \phi_Z^{(2)}Z_2 + \varepsilon_{X,3}^{(1)}, \text{ for } \varepsilon_{X,3}^{(2)} \sim N(0, 1);$$

(2) Generate  $Y_2$  using Equation (39)

$$Y_2 = \left( X_3 - \mu_X^{(1)} - \phi_X^{(1)}X_2 - \varepsilon_{X,3}^{(1)} \right) / \phi_Y^{(1)}, \text{ for } \varepsilon_{X,3}^{(1)} \sim N(0, 1) \text{ and } \varepsilon_{X,3}^{(1)} \perp \varepsilon_{X,3}^{(2)};$$

(3) Generate  $Z_{1,3}$  using Equation (41)

$$Z_{1,3} = \mu_Z - \phi_Y^{(2)} / \phi_Z^{(2)} Y_2 + \phi_Z Z_2 + \varepsilon_{Z1,3}, \text{ for } \varepsilon_{Z1,3} \sim N(0, 1) \text{ and } \varepsilon_{Z1,3} \perp \left( \varepsilon_{X,3}^{(1)}, \varepsilon_{X,3}^{(2)} \right);$$

(4) Generate the rest of  $Z_{i,t}$ , for  $i = 2, \dots, N$ , where  $Z_{t,i}$  are mutually independent, using

$$Z_{i,t} = \varepsilon_{Zi,t} \sim N(0, 1);$$

(5) Repeat steps (1)-(4)  $T + 500$  times and discard the first 500 observations to eliminate the effects of initial values.

To examine the performance of tests in Theorem 1, Table 1 summarizes the DGPs, with a direct causality from  $Y$  to  $X$ , that we use in our simulations and provides four different sets of parameters that we choose to indicate various degrees of causality from  $Y$  to  $X$  and different degrees of dependence in  $Z$ .

It is straightforward to notice that DGP1 and DGP3 are exhibiting a relatively weaker extent of direct causality from  $Y$  to  $X$  compared to DGP2 and DGP4, in terms of the coefficients in front of

Table 2: Data-generating processes: Indirect causality transmitted by one auxiliary variable

DGPs	Coefficients			
	Constants	$X$	$Y$	$Z$
DGP5	$\mu_X^{(1)} = \mu_X^{(2)} = \mu_Z = 0.5$	$\phi_X^{(1)} = \phi_X^{(2)} = 0.2$	$\phi_Y^{(1)} = \phi_Y^{(2)} = 0.3$	$\phi_Z^{(2)} = \phi_Z = 0.4$
DGP6	$\mu_X^{(1)} = \mu_X^{(2)} = \mu_Z = 0.5$	$\phi_X^{(1)} = \phi_X^{(2)} = 0.2$	$\phi_Y^{(1)} = \phi_Y^{(2)} = 0.7$	$\phi_Z^{(2)} = \phi_Z = 0.8$
DGP7	$\mu_X^{(1)} = \mu_X^{(2)} = \mu_Z = 0.5$	$\phi_X^{(1)} = \phi_X^{(2)} = 0.3$	$\phi_Y^{(1)} = \phi_Y^{(2)} = 0.2$	$\phi_Z^{(2)} = \phi_Z = 0.4$
DGP8	$\mu_X^{(1)} = \mu_X^{(2)} = \mu_Z = 0.5$	$\phi_X^{(1)} = \phi_X^{(2)} = 0.7$	$\phi_Y^{(1)} = \phi_Y^{(2)} = 0.2$	$\phi_Z^{(2)} = \phi_Z = 0.8$

**Note:** This table summarizes the DGPs with an indirect causality from  $Y$  to  $X$ , that we consider in the simulation study to examine the performance of the tests in Theorem 1 for testing the nulls outlined in (18), (21) and (22). The coefficients in this table are the coefficients of regression equations in (39)-(41). The error terms  $\varepsilon_{it}$  for  $i = 1, 2, 3$  and  $\varepsilon_{j+2,t}$  for  $j = 2, \dots, N$  are  $N + 2$  mutually independent standard normal random variables, where  $N$  can be large indicating the richness of the data environment.

$Y_{t-1}$ ; i.e. 0.3 versus 0.7. For the direct causality set-up in Table 1, the auxiliary variables  $Z$ s give no extra useful information for predicting  $X$  and we have used two different types of  $Z$ s, which can be either i.i.d. or AR(1) processes.

Furthermore, we consider the additional DGPs in Table 2 which correspond to equations (39)-(41) with four different sets of parameters that represent different scenarios of indirect causality. Since  $Z_1$  is present in DGP5 to DGP8, different parameter values indicate different degrees of indirect causality from  $Y$  to  $X$ . For instance, it seems at first sight that there is a high direct causality from  $Y$  to  $X$  in DGP6 given by  $\phi_Y^{(1)} = 0.7$ . However, this causality disappears once controlling the effect of  $Z$ , hence following into the context of indirect causality.

In the simulations, three different sample sizes are studied:  $T = 100, 200,$  and  $400$ . In addition,  $N$  is chosen to be varying according to the time periods  $T$ . In particular, for each DGP in tables (1) and (2), three different values of  $N$  are considered to examine the effect of data richness on the performance of the tests. The nominal size 5% is used and results for other significance levels are similar. Finally, all the results are based on 2000 replications.

Simulation results using DGP1 to DGP8 are reported in Tables 3-5. Observe that for testing condition (ii) of Definition (1), DGP1 to DGP4 are used to investigate the power performance of the tests, and DGP5 to DGP8 are used to examine their size properties [see Table 3 for details]. However, it is important to notice that the roles of null and alternative hypotheses are reversed when testing condition (iii) of Definition (1) using the regression equations (19) and (20) [see Tables 4 and

Table 3: Empirical rejection rates of the proposed test in regression (17) based on one  $Z$

DGPs								
	DGP1	DGP2	DGP3	DGP4	DGP5	DGP6	DGP7	DGP8
$T = 100$								
$N = 100$	32.9	68.2	34.3	66.3	5.4	5.4	5.0	5.5
$N = 200$	30.9	65.8	31.9	67.2	4.9	5.2	5.2	5.2
$N = 400$	32.4	65.4	29.3	65.2	5.0	5.0	5.3	8.0
$T = 200$								
$N = 200$	56.3	93.9	57.6	91.8	4.4	6.0	4.8	5.1
$N = 400$	55.9	91.7	57.9	93.1	5.0	5.7	4.6	6.3
$N = 600$	55.1	93.2	57.5	93.1	5.0	6.5	5.0	7.5
$T = 400$								
$N = 400$	87.4	99.9	85.8	99.9	7.3	5.2	6.1	6.9
$N = 600$	86.3	99.7	84.9	99.9	5.0	5.0	4.5	7.3
$N = 800$	86.4	99.9	85.4	99.9	5.3	4.6	4.6	8.6

**Note:** This table reports the empirical size and power of the test stated in Theorem 1 for testing condition (ii) of Definition 1 of indirect causality from  $Y$  to  $X$  at  $\alpha = 5\%$  significance level in regression (17). The number of simulations is equal to 2000 replications. The indirect causality is transmitted by one auxiliary variable.

Table 4: Empirical rejection rates of the proposed test in regression (19) based on one  $Z$

DGPs								
	DGP1	DGP2	DGP3	DGP4	DGP5	DGP6	DGP7	DGP8
$T = 100$								
$N = 100$	5.3	5.8	5.0	5.3	100.0	100.0	100.0	100.0
$N = 200$	4.6	6.0	4.8	5.2	100.0	100.0	100.0	100.0
$N = 400$	4.7	5.2	5.0	5.5	100.0	100.0	100.0	100.0
$T = 200$								
$N = 200$	5.6	4.7	5.1	5.4	100.0	100.0	100.0	100.0
$N = 400$	4.8	4.7	4.6	4.7	100.0	100.0	100.0	100.0
$N = 600$	5.9	4.9	5.6	5.4	100.0	100.0	100.0	100.0
$T = 400$								
$N = 400$	4.5	4.4	5.0	5.0	100.0	100.0	100.0	100.0
$N = 600$	4.7	4.9	6.1	4.9	100.0	100.0	100.0	100.0
$N = 800$	5.2	5.6	5.1	4.9	100.0	100.0	100.0	100.0

**Note:** This table reports the empirical size and power of the test stated in Theorem 2 for testing the null hypothesis (21) for the condition (iii) of Definition 1 of indirect causality from  $Y$  to  $X$  at  $\alpha = 5\%$  significance level in regression (19). The number of simulations is equal to 2000 replications. The indirect causality is transmitted by one auxiliary variable.

Table 5: Empirical rejection rates of the proposed test in regression (20) based on one  $Z$

DGP <sub>s</sub>								
	DGP1	DGP2	DGP3	DGP4	DGP5	DGP6	DGP7	DGP8
$T = 100$								
$N = 100$	5.2	5.5	4.5	4.8	48.9	99.7	38.8	41.6
$N = 200$	4.6	4.8	5.0	4.8	47.6	99.7	37.3	41.2
$N = 400$	5.5	5.4	5.9	4.8	47.5	98.8	37.4	39.2
$T = 200$								
$N = 200$	5.5	4.4	6.1	4.5	79.2	100.0	65.4	74.5
$N = 400$	5.0	5.8	5.5	5.7	79.4	100.0	65.9	74.1
$N = 600$	5.1	4.6	5.2	4.8	76.7	100.0	65.6	72.5
$T = 400$								
$N = 400$	5.0	4.3	5.3	4.4	97.5	100.0	92.7	96.6
$N = 600$	5.0	5.0	4.8	4.7	97.7	100.0	92.7	97.1
$N = 800$	5.2	4.4	4.9	5.0	96.7	100.0	92.5	96.4

**Note:** This table reports the empirical size and power of the test stated in Theorem 2 for testing the null hypothesis (22) for the condition (iii) of Definition 1 of indirect causality from  $Y$  to  $X$  at  $\alpha = 5\%$  significance level in regression (20). The number of simulations is equal to 2000. The indirect causality is transmitted by one auxiliary variable.



5 for details]. The results show that the empirical sizes are accurate for different DGPs with various sample sizes and numbers of  $Z$ s. Furthermore, it seems that the test of non-causality from  $Y$  to  $Z$  is substantially more powerful than the test of non-causality from  $Z$  to  $X$ . Finally, we find that the number of auxiliary variables  $Z$ s does not seem to affect greatly the performance of the tests given various sample sizes  $T$ .

To further illustrate the performance of the proposed tests, we consider the interesting case where many auxiliary variables [10 variables in our simulation] are transmitting the indirect causality from  $Y$  to  $X$ . In particular, we shall consider the following processes:

$$X_t = \mu_X^{(1)} + \phi_X^{(1)} X_{t-1} + \phi_Y^{(1)} Y_{t-1} + \varepsilon_{X,t}^{(1)}, \text{ with } \varepsilon_{X,t}^{(1)} \sim N(0, 1), \quad (42)$$

$$X_t = \mu_X^{(2)} + \phi_X^{(2)} X_{t-1} + \phi_Y^{(2)} Y_{t-2} + \sum_{j=1}^{10} \phi_{Z_j}^{(2)} Z_{j,t-1} + \varepsilon_{X,t}^{(2)}, \text{ with } \varepsilon_{X,t}^{(2)} \sim N(0, 1), \quad (43)$$

$$Z_{j,t} = \mu_{Z_j} - \frac{\phi_Y^{(2)}}{\phi_Z^{(2)}} Y_{t-1} + \phi_{Z_j} Z_{j,t-1} + \varepsilon_{Z_j,t}, \text{ with } \varepsilon_{Z_j,t} \sim N(0, 1), j = 1, \dots, 10, \text{ and } \phi_Z^{(2)} \neq 0, \quad (44)$$

$$Z_{i,t} = \varepsilon_{Z_{i,t}} \sim N(0, 1), \text{ for } i = 11, \dots, N, \text{ with } \varepsilon_{Z_{i,t}} \text{ mutually independent,}$$

where  $\phi_Z^{(2)} := \sum_{j=1}^{10} \phi_{Z_j}^{(2)}$ . The latter condition guarantees that  $Y$  does not cause  $X$  in the presence of  $Z$ , which satisfies condition (ii) of Definition (1). We now consider the five DGPs in Table 6 and follow the steps described above to simulate the data on  $X$ ,  $Y$ , and  $Z$ .

The DGPs in Table 6 are different from those in Tables 1 and 2. In particular, the values of parameters in the former are smaller, which indicates that the causal links are much weaker when we consider ten auxiliary variables. The use of smaller values is to ensure that the processes under consideration are stationary. Thus, the power of tests of indirect causality will be low when we use ten auxiliary variables instead of one.

Tables 7 to 9 report the empirical size and power of the tests of conditions (ii) and (iii) of an indirect causality from  $Y$  to  $X$  transmitted by 10 auxiliary variables. Before discussing the results, recall that DGPs 11 to 13 were used to assess the empirical size of test of condition (ii) [Theorem 1], whereas DGPs 9 and 10 were used to assess the empirical size of test of condition (iii) [Theorem 2]. The results show that the proposed tests control very well the size whatever the sample size  $T$ , the number of  $Z$ s, and the DGP under consideration. Regarding the power, on one hand we find that the test of condition (ii) in Theorem 1 has low power, but as expected it improves with the sample size. On the other hand, the power of test of condition (iii) in Theorem 2 is high and reaches one even when the sample size is small. Finally, our results indicate that both empirical size and power are quite stable when  $N$  changes.

Table 6: Data-generating processes: Direct causality and indirect causality transmitted by many auxiliary variables

DGPs	Variables of Interest		
	$X_t =$	$Y_t =$	$Z_t$
Direct Causality			
DGP9	$0.5 + 0.5X_{t-1} + 0.1Y_{t-1} + \varepsilon_{1t}$	$0.5 + 0.5Y_{t-1} + \varepsilon_{2t}$	$Z_{jt} = \varepsilon_{j+2,t}$ , for $j = 1, \dots, 10$ , $Z_{jt} = \varepsilon_{j+2,t}$ , for $j = 11, \dots, N$
DGP10	$0.5 + 0.5X_{t-1} + 0.1Y_{t-1} + \varepsilon_{1t}$	$0.5 + 0.5Y_{t-1} + \varepsilon_{2t}$	$Z_{jt} = 0.5 + 0.5Z_{j,t-1} + \varepsilon_{j+2,t}$ , for $j = 1, \dots, 10$ , $Z_{jt} = \varepsilon_{j+2,t}$ , for $j = 11, \dots, N$
Indirect Causality			
DGP11	$\phi_X^{(1)} = \phi_X^{(2)} = 0.1$	$\phi_Y^{(1)} = \phi_Y^{(2)} = 0.05$	$\phi_{Z_1}^{(2)} = \dots = \phi_{Z_{10}}^{(2)} = 0.01$ , $\phi_{Z_1} = \dots = \phi_{Z_{10}} = 0.2$
DGP12	$\phi_X^{(1)} = \phi_X^{(2)} = 0.2$	$\phi_Y^{(1)} = \phi_Y^{(2)} = 0.3$	$\phi_{Z_1}^{(2)} = \dots = \phi_{Z_{10}}^{(2)} = 0.05$ , $\phi_{Z_1} = \dots = \phi_{Z_{10}} = 0.4$
DGP13	$\phi_X^{(1)} = \phi_X^{(2)} = 0.2$	$\phi_Y^{(1)} = \phi_Y^{(2)} = 0.7$	$\phi_{Z_1}^{(2)} = \dots = \phi_{Z_{10}}^{(2)} = 0.1$ , $\phi_{Z_1} = \dots = \phi_{Z_{10}} = 0.8$

**Note:** This table summarizes the DGPs, with direct and indirect causalities from  $Y$  to  $X$ , that we consider in the simulation study to examine the performance of the tests for the hypothesis testing problems outlined in (18), (21) and (22). The error terms  $\varepsilon_{it}$  for  $j = 1, \dots, N + 2$  are mutually independent standard normal random variables, where  $N$  can be large indicating the richness of the data environment. The constant terms  $\mu_X^{(1)}$ ,  $\mu_X^{(2)}$ ,  $\mu_{Z_1}, \dots, \mu_{Z_{10}}$  in the equations (42), (43) and (44) for DGP11-DGP13 are all equal to  $0.5$ :  $\mu_X^{(1)} = \mu_X^{(2)} = \mu_{Z_1} = \dots = \mu_{Z_{10}} = 0.5$ .

Table 7: Empirical rejection rates of the proposed test in regression (17) based on ten  $Z$ s

DGPs					
	DGP9	DGP10	DGP11	DGP12	DGP13
$T = 100$					
$N = 100$	7.0	7.9	4.2	5.0	4.8
$N = 200$	9.0	9.7	5.6	4.8	5.8
$N = 400$	9.3	9.6	4.3	5.1	4.4
$T = 200$					
$N = 200$	12.0	13.3	4.4	4.8	4.8
$N = 400$	11.9	12.0	5.3	5.0	5.0
$N = 600$	12.7	14.0	5.2	4.3	5.0
$T = 400$					
$N = 400$	21.1	20.4	4.8	5.7	4.9
$N = 600$	20.7	19.2	5.0	5.4	4.6
$N = 800$	19.4	21.2	5.3	4.8	5.3

**Note:** This table reports the empirical size and power of the test stated in Theorem 1 for testing condition (ii) of Definition 1 of indirect causality from  $Y$  to  $X$  at  $\alpha = 5\%$  significance level in regression (17). The number of simulations is equal to 2000. The indirect causality is transmitted by ten auxiliary variables.

Table 8: Empirical rejection rates of the proposed test in regression (19) based on ten  $Z$ s

DGPs					
	DGP9	DGP10	DGP11	DGP12	DGP13
$T = 100$					
$N = 100$	5.9	5.5	100.0	100.0	100.0
$N = 200$	5.9	6.6	100.0	100.0	100.0
$N = 400$	5.6	5.7	100.0	100.0	100.0
$T = 200$					
$N = 200$	4.5	5.4	100.0	100.0	100.0
$N = 400$	6.2	4.6	100.0	100.0	100.0
$N = 600$	5.8	5.2	100.0	100.0	100.0
$T = 400$					
$N = 400$	4.6	4.9	100.0	100.0	100.0
$N = 600$	5.2	5.0	100.0	100.0	100.0
$N = 800$	5.2	6.7	100.0	100.0	100.0

**Note:** This table reports the empirical size and power of the test stated in Theorem 2 for testing the null hypothesis (21) of the condition (iii) of Definition 1 of indirect causality from  $Y$  to  $X$  at  $\alpha = 5\%$  significance level in regression (19). The number of simulations is equal to 2000. The indirect causality is transmitted by ten auxiliary variables.

Table 9: Empirical rejection rates of the proposed test in regression (20) based on ten  $Z$ s

DGP					
	DGP9	DGP10	DGP11	DGP12	DGP13
$T = 100$					
$N = 100$	5.3	5.6	8.0	32.7	99.5
$N = 200$	4.5	5.3	7.9	33.9	99.3
$N = 400$	4.9	5.1	8.6	31.8	99.0
$T = 200$					
$N = 200$	5.4	4.6	10.5	57.4	100
$N = 400$	5.5	5.9	10.2	58.6	100.0
$N = 600$	4.9	5.0	9.1	57.2	100.0
$T = 400$					
$N = 400$	4.9	4.3	14.4	85.4	100.0
$N = 600$	5.2	5.6	15.4	85.0	100.0
$N = 800$	4.7	4.9	14.8	85.4	100.0

**Note:** This table reports the empirical size and power of the test stated in Theorem 2 for testing the null hypothesis (22) of the condition (iii) of Definition 1 of indirect causality from  $Y$  to  $X$  at  $\alpha = 5\%$  significance level in regression (20). The number of simulations is equal to 2000. The indirect causality is transmitted by ten auxiliary variables.

## 7.2 Spurious causality

We now describe the DGPs that we use to evaluate the performance of tests of spurious causality of type I in Theorem 3. As for indirect causality, we first consider a set of DGPs in which the spurious causality of type I is induced by one and only one auxiliary variable  $Z_1$ , among a total of  $N$  variables  $\{Z_1, \dots, Z_N\}$ . In particular, we consider the following processes:

$$X_t = \mu_X^{(1)} + \phi_X^{(1)} X_{t-1} + \varepsilon_{X,t}^{(1)}, \text{ with } \varepsilon_{X,t}^{(1)} \sim N(0, 1) \quad (45)$$

$$X_t = \mu_X^{(2)} + \phi_X^{(2)} X_{t-1} + \phi_Z^{(2)} Z_{1,t-1} + \varepsilon_{X,t}^{(2)}, \text{ with } \varepsilon_{X,t}^{(2)} \sim N(0, 1) \quad (46)$$

$$Z_{1,t} = \mu_Z + \phi_Z^{(3)} Z_{1,t-1} + \phi_Y^{(3)} Y_{t-1} + \phi_X^{(3)} X_{t-1} + \varepsilon_{Z_1,t}, \text{ with } \varepsilon_{Z_1,t} \sim N(0, 1), \quad (47)$$

$$Z_{i,t} = \varepsilon_{Z_i,t} \sim N(0, 1), \text{ for } i = 2, \dots, N, \text{ with } \varepsilon_{Z_i,t} \text{ mutually independent,}$$

where the error terms  $\varepsilon_{X,t}^{(1)}$ ,  $\varepsilon_{X,t}^{(2)}$ , and  $\varepsilon_{Z_i,t}$ , for  $i = 1, \dots, N$ , are assumed to be mutually independent and  $X$  has to satisfy both equations (45) and (46). Numerical values for the coefficients  $\mu_X^{(1)}$ ,  $\mu_X^{(2)}$ ,  $\mu_Z$ ,  $\phi_X^{(1)}$ ,  $\phi_X^{(2)}$ ,  $\phi_X^{(3)}$ ,  $\phi_Z^{(2)}$ ,  $\phi_Z^{(3)}$ ,  $\phi_Y^{(3)}$  will be specified later. The functional forms of DGPs of  $X$  and  $Z$  are selected such that there is a spurious causality of type I from  $Y$  to  $X$ . According to Definition (2), spurious causality of type I occurs if: (i)  $Y$  Granger causes  $X$  with respect to the information set  $I(t) - I_Y(t)$ ; (ii)  $Y$  does not Granger cause  $X$  with respect to the information set  $I_X(t)$ ; and (iii)  $Y$  Granger causes  $Z$  and  $Z$  Granger causes  $X$  with respect to the information sets  $I(t) - I_Y(t)$  and  $I(t) - I_Z(t)$ , respectively. Thus, conditions (i) and (iii) will be satisfied if we choose the coefficients  $\phi_Z^{(2)}$  and  $\phi_Y^{(3)}$  in equations (46) and (47) to be different from zero. Furthermore, condition (ii) is also satisfied according to Equation (45). Now to help see why condition (i) is satisfied, observe that by combining equations (46) and (47) we obtain

$$X_t = \left( \mu_X^{(2)} + \phi_Z^{(2)} \mu_Z \right) + \phi_X^{(2)} X_{t-1} + \phi_Z^{(2)} \phi_X^{(3)} X_{t-2} + \phi_Z^{(2)} \phi_Y^{(3)} Y_{t-2} + \phi_Z^{(2)} \phi_Z^{(3)} Z_{1,t-2} + \phi_Z^{(2)} \varepsilon_{Z_1,t-1} + \varepsilon_{X,t}^{(2)},$$

where  $Y$  does cause  $X$  in the presence of  $Z$ .

Depending on whether or not  $\phi_Z^{(2)}$  and  $\phi_Y^{(3)}$  are taken to be equal to zero, the following steps can be performed to simulate a sample of  $T$  observations on  $X$ ,  $Y$  and  $Z$  under the absence/presence of spurious causality of type I from  $Y$  to  $X$  induced by the auxiliary variable  $Z_1$ :

(1) Choose the initial value  $X_1$  and generate  $X_2$  and  $X_3$  using Equation (45)

$$\begin{aligned} X_2 &= \mu_X^{(1)} + \phi_X^{(1)} X_1 + \varepsilon_{X,2}^{(1)}, \text{ for } \varepsilon_{X,2}^{(1)} \sim N(0, 1); \\ X_3 &= \mu_X^{(1)} + \phi_X^{(1)} X_2 + \varepsilon_{X,3}^{(1)}, \text{ for } \varepsilon_{X,3}^{(1)} \sim N(0, 1) \text{ and } \varepsilon_{X,2}^{(1)} \perp \varepsilon_{X,3}^{(1)}; \end{aligned}$$

Table 10: Data-generating processes: Non-causality cases

DGPs	Variables of Interest		
	$X_t =$	$Y_t =$	$Z_{jt}$ for $j = 1, \dots, N$
DGP14	$0.5 + 0.8X_{t-1} + \varepsilon_{1t}$	$0.5 + 0.8Y_{t-1} + \varepsilon_{2t}$	$Z_{1t} = \varepsilon_{3t}, Z_{jt} = \varepsilon_{j+2,t}$
DGP15	$0.5 + 0.2X_{t-1} + \varepsilon_{1t}$	$0.5 + 0.2Y_{t-1} + \varepsilon_{2t}$	$Z_{1t} = \varepsilon_{3t}, Z_{jt} = \varepsilon_{j+2,t}$
DGP16	$0.5 + 0.8X_{t-1} + \varepsilon_{1t}$	$0.5 + 0.8Y_{t-1} + \varepsilon_{2t}$	$Z_{1t} = 0.5 + 0.8Z_{1,t-1} + \varepsilon_{3t}, Z_{jt} = \varepsilon_{j+2,t}$
DGP17	$0.5 + 0.2X_{t-1} + \varepsilon_{1t}$	$0.5 + 0.2Y_{t-1} + \varepsilon_{2t}$	$Z_{1t} = 0.5 + 0.2Z_{1,t-1} + \varepsilon_{3t}, Z_{jt} = \varepsilon_{j+2,t}$

**Note:** This table summarizes the DGPs, with no spurious causality of type I from  $Y$  to  $X$ , that we consider in the simulation study to examine the size of the tests for the hypothesis testing problems outlined in (24), (28), and (29). The error terms  $\varepsilon_{it}$  for  $j = 1, \dots, N + 2$  are mutually independent standard normal random variables, where  $N$  can be large indicating the richness of the data environment.

(2) Generate  $Z_{1,1}$  and  $Z_{1,2}$  using Equation (46)

$$Z_{1,1} = \left( X_2 - \mu_X^{(2)} - \phi_X^{(2)} X_1 - \varepsilon_{X,2}^{(2)} \right) / \phi_Z^{(2)}, \text{ for } \phi_Z^{(2)} \neq 0, \varepsilon_{X,2}^{(2)} \sim N(0, 1), \text{ and } \varepsilon_{X,2}^{(2)} \perp \left( \varepsilon_{X,2}^{(1)}, \varepsilon_{X,3}^{(1)} \right)$$

$$Z_{1,2} = X_3 - \mu_X^{(2)} - \phi_X^{(2)} X_2 - \varepsilon_{X,3}^{(2)} / \phi_Z^{(2)}, \text{ for } \phi_Z^{(2)} \neq 0, \varepsilon_{X,3}^{(2)} \sim N(0, 1), \text{ and } \varepsilon_{X,3}^{(2)} \perp \left( \varepsilon_{X,2}^{(2)}, \varepsilon_{X,2}^{(1)}, \varepsilon_{X,3}^{(1)} \right);$$

(3) Generate  $Y_1$  using Equation (47)

$$Y_1 = \left( Z_{1,2} - \mu_Z - \phi_Z^{(3)} Z_{1,1} - \phi_X^{(3)} X_1 - \varepsilon_{Z,1,2} \right) / \phi_Y^{(3)}, \text{ for } \varepsilon_{Z,1,2} \sim N(0, 1) \text{ and } \varepsilon_{Z,1,2} \perp \left( \varepsilon_{X,2}^{(2)}, \varepsilon_{X,2}^{(1)}, \varepsilon_{X,3}^{(1)} \right);$$

(4) Generate the rest of  $Z_{i,t}$ , for  $i = 2, \dots, N$ , where  $Z_{t,i}$  are mutually independent, using

$$Z_{i,t} = \varepsilon_{Z,i,t} \sim N(0, 1);$$

(5) Repeat steps (1)-(4)  $T + 500$  times and discard the first 500 observations to eliminate the effects of the selection of initial values.

To examine the empirical size of tests in Theorem 3, Table 10 summarizes the DGPs, when there is no causality from  $Y$  to  $X$ , that we use in our simulations. Regarding the assessment of the empirical power, we consider the DGPs in Table 11 that correspond to equations (45), (46), and (47) with four different sets of parameters that represent different scenarios of spurious causality of type I. As in Section 7.1, three different sample sizes are studied;  $T = 100, 200$ , and  $400$ , and  $N$  is chosen to be varying according to the sample sizes  $T$ . The nominal size 5% is studied and results for other significance levels are similar. Finally, all the results are based on 2000 replications.

Tables (12) to (14) report the empirical size and power of tests of conditions of spurious causality of Type I under the DGPs in tables (10) and (11). On the one hand, the results, using DGPs 14 to 17,

Table 11: Data-generating processes: Spurious causality of type I

DGPs	Coefficients			
	Constants	$X$	$Y$	$Z$
DGP18	$\mu_X^{(1)} = \mu_X^{(2)} = \mu_Z = 0.5$	$\phi_X^{(1)} = \phi_X^{(2)} = \phi_X^{(3)} = 0.1$	$\phi_Y^{(3)} = 0.1$	$\phi_Z^{(2)} = \phi_Z^{(3)} = 0.1$
DGP19	$\mu_X^{(1)} = \mu_X^{(2)} = \mu_Z = 0.5$	$\phi_X^{(1)} = \phi_X^{(2)} = \phi_X^{(3)} = 0.2$	$\phi_Y^{(3)} = 0.2$	$\phi_Z^{(2)} = \phi_Z^{(3)} = 0.2$
DGP20	$\mu_X^{(1)} = \mu_X^{(2)} = \mu_Z = 0.5$	$\phi_X^{(1)} = \phi_X^{(2)} = \phi_X^{(3)} = 0.3$	$\phi_Y^{(3)} = 0.3$	$\phi_Z^{(2)} = \phi_Z^{(3)} = 0.3$
DGP21	$\mu_X^{(1)} = \mu_X^{(2)} = \mu_Z = 0.5$	$\phi_X^{(1)} = \phi_X^{(2)} = \phi_X^{(3)} = 0.4$	$\phi_Y^{(3)} = 0.4$	$\phi_Z^{(2)} = \phi_Z^{(3)} = 0.4$

**Note:** This table summarizes the DGPs, with a spurious causality of type I from  $Y$  to  $X$ , that we consider in the simulation study to examine the power of the tests for the hypothesis testing problems outlined in (24), (28), and (29). The error terms  $\varepsilon_{it}$  for  $j = 1, \dots, N + 2$  are mutually independent standard normal random variables, where  $N$  can be large indicating the richness of the data environment.

show that the proposed tests control the size whatever the sample size and the number of auxiliary variables. The size control is achieved by all the tests in Theorem 3 and under all the DGPs, except DGP 16 for which the empirical size is slightly higher than the significance level of 5%. On the other hand, the empirical power of the tests reaches one for all DGPs [DGP18 to DGP21], even when the sample size is small and whatever the number of auxiliary variables. Finally, the case of multiple auxiliary variable  $Z$ s is omitted for sake of brevity.

## 8 Empirical application

We use the tests proposed in the above sections to test for the presence of an indirect causality from credit/money to real activity. Studying the interaction between real activity (income) and monetary policy measures (money) and credit is of great importance to economists, see e.g. Friedman (1981).

As pointed out by Fackler (1985), in studying the relationship between money and income, empirical evidence suggests that important information may be lost by ignoring some variables such as the one that comes from credit market. He argued that empirical results on examining money-income causal relationship might differ depending on the information set one has at hands. His empirical analysis shows that the results based on bivariate causality analysis are misleading and often overturned when one extends the information set and includes other key variables such as the ones related to the credit. In particular, he found that interest rates play the role of an auxiliary variable that transmits the causality between financial and real sectors. In other words, money/credit does not



Table 12: Empirical rejection rates of the proposed test in regression (23) based on one  $Z$

DGPs								
	DGP14	DGP15	DGP16	DGP17	DGP18	DGP19	DGP20	DGP21
$T = 100$								
$N = 100$	6.6	4.7	5.9	5.3	100.0	100.0	100.0	100.0
$N = 200$	6.6	5.1	7.2	5.8	100.0	100.0	100.0	100.0
$N = 400$	6.7	5.7	6.5	5.4	100.0	100.0	100.0	100.0
$T = 200$								
$N = 200$	6.4	5.0	6.9	5.6	100.0	100.0	100.0	100.0
$N = 400$	5.7	5.1	6.6	5.2	100.0	100.0	100.0	100.0
$N = 600$	6.2	5.5	6.9	5.4	100.0	100.0	100.0	100.0
$T = 400$								
$N = 400$	5.2	5.2	5.8	5.3	100.0	100.0	100.0	100.0
$N = 600$	5.2	4.9	5.6	4.7	100.0	100.0	100.0	100.0
$N = 800$	5.6	5.5	5.3	4.6	100.0	100.0	100.0	100.0

**Note:** This table reports the empirical size and power of the test stated in Theorem 3 for testing condition (i) of Definition 2 of spurious causality of type I from  $Y$  to  $X$  at  $\alpha = 5\%$  significance level in regression (23). The number of simulations is equal to 2000.

Table 13: Empirical rejection rates of the proposed test in regression (26) based on one  $Z$

DGPs								
	DGP14	DGP15	DGP16	DGP17	DGP18	DGP19	DGP20	DGP21
$T = 100$								
$N = 100$	5.3	5.4	7.3	5.4	100.0	100.0	100.0	100.0
$N = 200$	5.3	4.8	8.2	5.1	100.0	100.0	100.0	100.0
$N = 400$	5.6	5.3	7.4	4.1	100.0	100.0	100.0	100.0
$T = 200$								
$N = 200$	4.6	5.0	7.6	5.1	100.0	100.0	100.0	100.0
$N = 400$	4.8	4.6	7.1	5.1	100.0	100.0	100.0	100.0
$N = 600$	5.5	4.7	6.7	5.2	100.0	100.0	100.0	100.0
$T = 400$								
$N = 400$	5.0	5.3	8.0	5.1	100.0	100.0	100.0	100.0
$N = 600$	6.2	3.9	8.0	5.4	100.0	100.0	100.0	100.0
$N = 800$	5.0	4.6	7.6	5.0	100.0	100.0	100.0	100.0

**Note:** This table reports the empirical size and power of the test stated in Theorem 3 for testing the null hypothesis (28) of the condition (iii) of Definition 2 of spurious causality of type I from  $Y$  to  $X$  at  $\alpha = 5\%$  significance level in regression (26). The number of simulations is equal to 2000.

Table 14: Empirical rejection rates of the proposed test in regression (27) based on one  $Z$

DGPs								
	DGP14	DGP15	DGP16	DGP17	DGP18	DGP19	DGP20	DGP21
$T = 100$								
$N = 100$	5.0	5.4	7.0	4.8	100.0	100.0	100.0	100.0
$N = 200$	4.2	4.6	5.7	5.7	100.0	100.0	100.0	100.0
$N = 400$	4.7	4.4	6.4	5.6	100.0	100.0	100.0	100.0
$T = 200$								
$N = 200$	4.2	5.7	5.1	5.0	100.0	100.0	100.0	100.0
$N = 400$	5.1	5.1	5.6	4.1	100.0	100.0	100.0	100.0
$N = 600$	4.7	5.0	4.8	4.9	100.0	100.0	100.0	100.0
$T = 400$								
$N = 400$	5.5	4.5	5.5	4.5	100.0	100.0	100.0	100.0
$N = 600$	4.8	5.0	4.5	4.7	100.0	100.0	100.0	100.0
$N = 800$	4.3	5.3	5.9	5.3	100.0	100.0	100.0	100.0

**Note:** This table reports the empirical size and power of the test stated in Theorem 3 for testing the null hypothesis (29) of the condition (iii) of Definition 2 of spurious causality of type I from  $Y$  to  $X$  at  $\alpha = 5\%$  significance level in regression (27). The number of simulations is equal to 2000.

directly influence real output; but it plays at most an indirect role in income determination.

To obtain his results, Fackler (1985) applied an ad hoc approach in which the auxiliary variables were predetermined or specified at the beginning of the analysis and not selected by any statistical method. Furthermore, his tests were run in the presence of only few variables, thus this excluded the hundred of economic variables that might play a role in income determination.

Our objective is to use the tests proposed in Section 5.1 to re-examine the existence of an indirect causality from money/credit to income using a recent dataset that contains more than 130 economic variables. In particular, we would like to confirm whether or not there is an indirect causality from money/credit to income. Thereafter, if this indirect causality exists, we would like to use the algorithm discussed in Section 6 to identify the auxiliary variable(s) that transmit this causality and compare them with those used in Fackler (1985). In this application income is measured by Industrial Production Index (IPI), money is measured by M1 Money Stock [hereafter M1SL using the FRED], and we consider two measures of credit: Commercial and Industrial Loans [hereafter BUSLOANS] and Securities in Bank Credit at All Commercial Banks [hereafter INVEST].

## 8.1 Data

We consider a big data set that consists of monthly observations on 135 economic variables from Federal Reserve Bank of St. Louis (FRED). The sample runs from January 1959 to May 2016 for a total of 689 observations. All the variables are reported in Tables 17-22 of Appendix B. In particular, we consider 8 groups of variables: (1) Output and income with 17 variables; (2) Labor market with 32 variables; (3) Housing with 10 variables; (4) Consumption, orders, and inventories with 14 variables; (5) Money and credit with 14 variables; (6) Interest and exchange rates with 22 variables; (7) Prices with 21 variables; and (8) Stock market with 5 variables. This big data mimic the coverage of datasets already used in the literature and it is updated in real-time through the FRED database. A detailed description of the dataset can be found in McCracken and Ng (2015).

## 8.2 Results

First of all, using Akaike information criterion, our results show that regressions (9), (17), (19), and (20) with 3 or 4 lagged terms are appropriate to test the conditions of a possible indirect causality from money/credit to income.

Table 15 reports the  $p$ -values for testing the conditions in Definition (1). On one hand, we find

Table 15: Testing Indirect Causality between Income and Money, Credit

Tested conditions		Causal variable ( $Y$ )		
		Money	BUSLOANS	INVEST
(i):	$Y \rightarrow X$	0.8655	0.0728	0.0717
(ii):	$Y \rightarrow X   f$	0.8098	0.5618	0.4222
(iii).a:	$Y \rightarrow f   X$	0.0195	0.0385	0.0011
(iii).b:	$f \rightarrow X   Y$	0.0012	0.0005	0.0009

**Note:** This table summarizes the results of testing conditions (i)-(iii) of Definition (1) of an indirect causality from  $Y$  =Money, Commercial and Industrial Loans [BUSLOANS] and Securities in Bank Credit at All Commercial Banks [INVEST] to  $X$  =Income.  $f$  represents the factor extracted from the big data set [see Section 4.1]. Money is measured by M1 Money Stock (M1SL) in Table 19 of Appendix B.

that there is no indirect causality from money to income, as the first condition [money Granger causes income without the presence of other variables] is not satisfied with a significantly high  $p$ -value. Consequently, this renders the subsequent testing procedure unnecessary, even though all the following conditions are satisfied. Hence, we conclude that money is not Granger indirectly causing income, which goes against the findings in Fackler (1985). In Table 15, we only include the results from one measure of money, i.e. M1SL, for illustration and as a matter of fact all the other money measures fail to pass the indirect causality tests and demonstrate quantitatively similar results.

On the other hand, the measures of credit [BUSLOANS and INVEST] fit to our testing paradigm well. In particular, the  $p$ -value for the first condition [BUSLOANS Granger causing income without the presence of other variables] is equal to 0.0728. However, once the auxiliary variable  $f$  [extracted from 135 economic variables] is included, BUSLOANS does not Granger cause income any more with a high  $p$ -value of 0.5618. Lastly, for the third condition, BUSLOANS appears to Granger cause the auxiliary variable  $f$  and  $f$  furthermore Granger causes income with  $p$ -values of 0.0385 and 0.0005, respectively. This leads us to believe that Commercial and Industrial Loans serves as an indirect source of income. For the credit measure INVEST, the same argument applies and the four  $p$ -values again help us to conclude that INVEST Granger causes income, but only in an indirect way. These results are in line with the findings in Fackler (1985).

Now that we have found that there is an indirect causality from credit to income, we next use the procedure outlined in Section 6 to identify the auxiliary variables that transmit this causality. The results are summarized in Table 16, where the first and second columns report the auxiliary variables

Table 16: Identification of the auxiliary variables responsible for the transmission of indirect causality from Credit, Investment to Income

Auxiliary variables ( $Z$ )	
Indirect causality from BUSLOANS to Income	Indirect causality from INVEST to Income
All Employees: Mining and Logging: Mining	IP: Nondurable Materials
Avg Weekly Overtime Hours : Manufacturing	Housing Starts: Total New Privately Owned
New Private Housing Permits (SAAR)	Housing Starts, West
New Private Housing Permits, South (SAAR)	New Private Housing Permits (SAAR)
New Private Housing Permits, West (SAAR)	New Private Housing Permits, Midwest (SAAR)
Total Business: Inventories to Sales Ratio	New Private Housing Permits, South (SAAR)
Effective Federal Funds Rate	New Private Housing Permits, West (SAAR)
3-Month AA Financial Commercial Paper Rate	Real Manu. and Trade Industries Sales
1-Year Treasury Rate	Total Business: Inventories to Sales Ratio
Moody's Seasoned Baa Corporate Bond Yield	5-Year Treasury C Minus FEDFUNDS
3-Month Treasury C Minus FEDFUNDS	10-Year Treasury C Minus FEDFUNDS
6-Month Treasury C Minus FEDFUNDS	Moody's Aaa Corporate Bond Minus FEDFUNDS
1-Year Treasury C Minus FEDFUNDS	Moody's Baa Corporate Bond Minus FEDFUNDS
5-Year Treasury C Minus FEDFUNDS	S&P's Composite Common Stock: Dividend Yield
10-Year Treasury C Minus FEDFUNDS	-
Moody's Aaa Corporate Bond Minus FEDFUNDS	-
Moody's Baa Corporate Bond Minus FEDFUNDS	-
S&P's Composite Common Stock: Dividend Yield	-

**Note:** This table summarizes the results of identifying the auxiliary variables responsible for the transmission of indirect causality from credit measured by Commercial and Industrial Loans and Securities in Bank Credit at All Commercial Banks to Income. The results are obtained using the statistical procedure described in Section 6 and based on the big data in Appendix B.

that transmit the indirect causality from Commercial and Industrial Loans and Securities in Bank Credit at All Commercial Banks to income, respectively. On one hand, we see that 18 auxiliary variables are responsible for the transmission of indirect causality from BUSLOANS to income. These variables belong to five groups: (i) Labor market; (ii) Housing; (iii) Consumption, orders, and inventories; (iv) Interest and exchange rates; and (v) Stock market. The variables from the other groups are found to be silent. We also find that most of the auxiliary variables [11 over a total of 18] belong to the group on interest and exchange rates. These variables are: Effective Federal Funds Rate, 3-Month AA Financial Commercial Paper Rate, 1-Year Treasury Rate, Moody's Seasoned Baa Corporate Bond Yield, 3-Month Treasury C Minus FEDFUNDS, 6-Month Treasury C Minus FEDFUNDS, 1-Year Treasury C Minus FEDFUNDS, 5-Year Treasury C Minus FEDFUNDS, 10-Year Treasury C Minus FEDFUNDS, Moody's Aaa Corporate Bond Minus FEDFUNDS, and Moody's Baa Corporate Bond Minus FEDFUNDS. Thus, it seems that the short and long-term interest rates are the main auxiliary variables that transmit the indirect causality from BUSLOANS to income, which is in line with the findings in Fackler (1985). Fackler (1985) wrote: *“What is presumably relevant for income determination, and especially for the investment component of income, is the long-term interest rate.”*

Regarding the indirect causality from INVEST to income, column 2 of Table 16 shows that 14 auxiliary variables are responsible for the transmission of this causality. These variables belong to five groups: (i) Output and income; (ii) Housing; (iii) Consumption, ordered inventories; (iv) Interest and exchange rates; and (v) Stock market. The dominant groups with highest numbers of auxiliary variables are Housing and Interest and exchange rates. Thus, in addition to the short and long-term interest rates, the housing sector is essential for transmitting the causality from credit to income, which is different from the findings in Fackler (1985).

## 9 Conclusion

We introduced a novel statistical procedure for testing indirect and spurious causal effects. Ignoring these effects might lead to wrong economic analysis, and consequently to inaccurate policy decisions. In practice, detecting indirect/spurious causality is a complicated task, since the pertinent auxiliary variables that transmit/induce the indirect/spurious causality are very often unknown. The availability of hundreds of economic variables makes this task even harder as it is generally infeasible to find the appropriate auxiliary variable(s) among all the available ones. In addition, including

hundreds of variables and their lags in a regression equation is technically difficult. We proposed new statistical procedures to test for the presence of an indirect/spurious causality using big data analysis. A diffusion index was included in the regression equation to represent all the variables that are available to practitioners. We derived the asymptotic distributions of the tests in the presence of an estimated index. Furthermore, we conducted a Monte Carlo simulation to evaluate the performance of the proposed statistical procedure. The results showed that our procedure is efficient for detecting indirect/spurious causality. Finally, we provided an empirical application where hundreds of variables are used to study a possible indirect causality from money/credit to income.

## References

- [1] Bai, J. (2003). “Inferential theory for factor models of large dimensions,” *Econometrica*, vol. 71, pp. 135–171.
- [2] Bai, J., Ng, S. (2002). “Determining the number of factors in approximate factor models,” *Econometrica*, vol. 70, pp. 191–221.
- [3] Bai, J., Ng, S. (2005). “Confidence intervals for diffusion index forecasts and inference for factor-augmented regressions,” *Working paper*, University of Michigan.
- [4] Bai, J., Ng, S. (2006). “Confidence intervals for diffusion index forecasts and inference for factor-augmented regressions,” *Econometrica*, vol. 74, pp. 1133–1150.
- [5] Bai, J., Ng, S. (2008). “Forecasting economic time series using targeted predictors,” *Journal of Econometrics*, vol. 146, pp. 304–317.
- [6] Dufour, J.M., Renault, E. (1998). “Short run and long run causality in time series: theory,” *Econometrica*, pp.1099–1125.
- [7] Fackler, J.S. (1985). “An empirical analysis of the markets for goods, money, and credit,” *Journal of Money, Credit and Banking*, vol. 17, pp.28–42.
- [8] Friedman, B.M. (1981). “The roles of money and credit in macroeconomic analysis,” *Mimeographed*.
- [9] Granger, C.W.J. (1969). “Investigating causal relations by econometric models and cross-spectral methods,” *Econometrica*, vol. 37, pp. 424–459.



- [10] Hsiao, C. (1982). “Autoregressive modeling and causal ordering of economic variables,” *Journal of Economic Dynamics and Control*, vol. 4, pp. 243–259.
- [11] Ludvigson, S., Ng, S. (2009). “Macro factors in bond risk premia,” *Review of Financial Studies*, vol. 22, p.p 5027–5067.
- [12] Lütkepohl, H. (1993). “Testing for causation between two variables in higher-dimensional VAR models,” In *Studies in Applied Econometrics*, pp. 75–91, *Physica-Verlag HD*.
- [13] McCracken, M.W., Ng, S. (2015). “FRED-MD: A monthly database for macroeconomic research,” Federal Reserve Bank of St. Louis Working Paper Series.
- [14] Stock, J.H., Watson, M.W. (2002). “Forecasting using principal components from a large number of predictors,” *Journal of the American Statistical Association*, vol. 97, pp. 1167–1179.
- [15] Wiener, N. (1956): “The theory of prediction,” in *E. F. Beckenback, ed., The Theory of Prediction*, McGraw-Hill, New York, chapter 8.

## A Appendix: Proofs

This appendix provides the proofs of Theorems 1 to 3 in the main text. We first introduce some notations, which are adapted from Bai and Ng (2006). Let  $\hat{V}$  be the  $(k \times k)$  diagonal matrix consisting of the  $k$  largest eigenvalues of  $WW'/(TN)$  and let  $H = \hat{V}^{-1}(\hat{f}'f/T)(\Lambda'\Lambda/N)$  be the rotation matrix, due to the fact that  $\hat{f}$  can only consistently estimate  $Hf$ , the space spanned by the true factor  $f$ . Let  $\Phi_0 = \text{diag}(I_{1+\bar{p}+\bar{q}}, V^{-1}Q\Sigma_\Lambda)$  being block diagonal, where  $V = \text{plim } \hat{V}$ ,  $Q = \text{plim } \hat{f}'f/T$  and  $\Sigma_\Lambda$  is defined in Assumption **A**.

Three auxiliary lemmas are first given below. The first one is due to Bai and Ng (2005).

**Lemma 1:** Take  $\hat{z}_t$  to be  $\hat{z}_{jt}$ , or  $\hat{z}_{jt}^{sii}$ , or  $\hat{z}_{jt}^{si}$  for any  $j = 1, 2, 3$ , which are defined in Section 5. Let  $z_t$  be the corresponding infeasible regressors and  $\bar{e}_{t+1}$  be any of the error terms in the corresponding regression. Let  $\delta_{NT}^2 = \min[N, T]$ . Then under Assumption **A**, we have: **(i)**  $\frac{1}{T} \sum_{t=1}^T \|\hat{f}_t - Hf_t\|^2 = O_p(\delta_{NT}^{-2})$ ; **(ii)**  $\frac{1}{T} \sum_{t=1}^T (\hat{f}_t - Hf_t)z_t' = O_p(\delta_{NT}^{-2})$ ; **(iii)**  $\frac{1}{T} \sum_{t=1}^T (\hat{f}_t - Hf_t)\hat{z}_t' = O_p(\delta_{NT}^{-2})$ ; and **(iv)**  $\frac{1}{T} \sum_{t=1}^T (\hat{f}_t - Hf_t)\bar{e}_{t+1}' = O_p(\delta_{NT}^{-2})$ .

**Lemma 2:** Let  $\hat{z}_t$  and  $z_t$  be the feasible and infeasible regressors defined in Lemma 1. Let  $\delta_{NT}^2 = \min[N, T]$ . Then under Assumption **A**, we have: **(i)**  $\frac{1}{T} \sum_{t=1}^T (\hat{f}_{t+1} - Hf_{t+1})z_t' = O_p(\delta_{NT}^{-2})$  and **(ii)**  $\frac{1}{T} \sum_{t=1}^T (\hat{f}_{t+1} - Hf_{t+1})\hat{z}_t' = O_p(\delta_{NT}^{-2})$ .

**Proof of Lemma 2:** They are similar to the proofs of results (ii) and (iii) in Lemma 1. ■

**Lemma 3:** Consider the infeasible regression (17) or (20). Let  $\tau$  be the parameter to be estimated and  $\hat{\tau}$  its ordinary least squares estimate obtained from a regression of  $X_{t+1}$  on the vector of regressors  $\hat{z}_t$ , with  $\hat{z}_t$  includes the intercept, lagged values of  $X_t, Y_t$ , and the estimated factor  $\hat{f}_t$  and its lags. Suppose Assumption **A** hold. If  $\sqrt{T}/N \rightarrow 0$ , then

$$\sqrt{T}(\hat{\tau} - \tau) \rightarrow_d N(0, \Sigma_\tau),$$

where  $\Sigma_\tau = \Phi_0'^{-1} \Sigma_{zz}^{-1} \Sigma_{zz,e} \Sigma_{zz}^{-1} \Phi_0^{-1}$ . Moreover, a consistent estimator of the variance covariance matrix  $\Sigma_\tau$  is given by

$$\hat{\Sigma}_\tau = \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_t \hat{z}_t' \right)^{-1} \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{e}_{t+1}^2 \hat{z}_t \hat{z}_t' \right) \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_t \hat{z}_t' \right)^{-1},$$

where  $\hat{e}_{t+1} = X_{t+1} - \hat{z}_t' \hat{\tau}$  are the least squares residuals.

**Proof of Lemma 3:** For Lemma 3, we will only prove the case for regression (17) as the proof for the limiting distribution of  $\hat{\tau}$  in regression (20) is identical and hence it is omitted. The proof is much similar to that of Theorem 1 in Bai and Ng (2006). Without loss of generality, assume  $\bar{p} = \bar{q} = \bar{h} = 1$ , and define  $z_t = (1, X_t, Y_t, f_t)'$  and  $\hat{z}_t = (1, X_t, Y_t, \hat{f}_t)'$  so that  $\tau = (\eta, \gamma, \lambda, \theta H^{-1})'$  are the parameters from the infeasible regression when  $f_t$  is observed.

From the infeasible regression (17), adding and subtracting terms, we obtain

$$\begin{aligned} X_{t+1} &= \eta + \gamma X_t + \lambda Y_t + \theta f_t + e_{t+1} = \eta + \gamma X_t + \lambda Y_t + \theta H^{-1} \hat{f}_t + e_{t+1} + \theta H^{-1} (H f_t - \hat{f}_t) \\ &= \hat{z}_t' \tau + e_{t+1} + \theta H^{-1} (H f_t - \hat{f}_t). \end{aligned}$$

In matrix notation,  $X = \hat{z} \tau + e + (f H' - \hat{f}) H^{-1} \theta$ , where  $X = (X_2, \dots, X_T)'$ ,  $\hat{z} = (\hat{z}_1, \dots, \hat{z}_{T-1})'$ ,  $e = (e_2, \dots, e_T)'$ ,  $f = (f_1, \dots, f_{T-1})'$  and  $\hat{f} = (\hat{f}_1, \dots, \hat{f}_{T-1})'$ . Therefore, the ordinary least squares estimator of  $\tau$  is given by

$$\hat{\tau} = (\hat{z}' \hat{z})^{-1} \hat{z}' X = \tau + (\hat{z}' \hat{z})^{-1} \hat{z}' e + (\hat{z}' \hat{z})^{-1} \hat{z}' (f H' - \hat{f}) H^{-1} \theta.$$

Thus,

$$\sqrt{T}(\hat{\tau} - \tau) = (T^{-1} \hat{z}' \hat{z})^{-1} T^{-1/2} \hat{z}' e + (T^{-1} \hat{z}' \hat{z})^{-1} [T^{-1/2} \hat{z}' (f H' - \hat{f})] H^{-1} \theta.$$

By the result (iii) of Lemma 1, the second term on the right-hand side of the above equation is  $O_p(T^{1/2}/\min(N, T)) = o_p(1)$  if  $\sqrt{T}/N \rightarrow 0$ . Define  $W_t = (1, X_t, Y_t)'$  so that  $T^{-1/2} \hat{z}' e = T^{-1/2} (e' W, e' \hat{f})'$ . Due to the fact that  $T^{-1/2} \hat{f}' e = T^{-1/2} H f' e + T^{-1/2} (\hat{f} - f H')' e$ , we have

$$T^{-1/2} \hat{f}' e = T^{-1/2} H f' e + o_p(1)$$

by the result (iv) of Lemma 1 and  $\sqrt{T}/N \rightarrow 0$ . Thus, we get that  $T^{-1/2}\hat{f}'e = T^{-1/2}(e'W, e'fH')' + o_p(1) = T^{-1/2}\Phi z'e + o_p(1)$ , with  $\Phi = \text{diag}(I, H)$  a block diagonal matrix. Therefore,

$$\sqrt{T}(\hat{\tau} - \tau) = (T^{-1}\hat{z}'\hat{z})^{-1}T^{-1/2}z'e + o_p(1) = (T^{-1}\hat{z}'\hat{z})^{-1}\Phi T^{-1/2}z'e + o_p(1).$$

Under standard assumptions, we have  $T^{-1/2}\sum_{t=1}^{T-1} z_t e_{t+1} \rightarrow_d N(0, \Sigma_{zz,e})$  with  $\Sigma_{zz,e} = \text{plim} T^{-1} \sum_{t=1}^{T-1} e_{t+1}^2 z_t z_t'$ . Therefore,  $\sqrt{T}(\hat{\tau} - \tau) \rightarrow_d N(0, \Sigma_\tau)$  with the asymptotic variance covariance matrix given by

$$\Sigma_\tau = \text{plim} \left( \frac{\hat{z}'\hat{z}}{T} \right)^{-1} \Phi \left( \frac{1}{T} \sum_{t=1}^{T-1} e_{t+1}^2 z_t z_t' \right) \Phi' \left( \frac{\hat{z}'\hat{z}}{T} \right)^{-1},$$

where  $\Phi = \text{diag}(I, H)$  is a block diagonal matrix with the probability limit  $\Phi_0$ . Following Bai and Ng (2006),  $\Sigma_\tau = \Phi_0^{-1} \Sigma_{zz}^{-1} \Sigma_{zz,e} \Sigma_{zz}^{-1} \Phi_0^{-1}$ .

In addition, by Bai and Ng (2006),  $\hat{\Sigma}_\tau$  is a consistent estimator for  $\Sigma_\tau$ . ■

**Proof of Theorem 1:** We focus on the case where  $\bar{p} = \bar{q} = \bar{h} = 1$ . Under heteroskedasticity, the corresponding  $F$ -statistic is defined by

$$F_T^{Ind,\lambda} = \left( \sqrt{T} R^{Ind,\lambda \hat{\tau}} \right) \left( R^{Ind,\lambda \hat{\Sigma}_\tau} R^{Ind,\lambda'} \right)^{-1} \left( \sqrt{T} R^{Ind,\lambda \hat{\tau}} \right)',$$

where  $R^{Ind,\lambda} = (0, 0, 1, 0)$ . Since  $R^{Ind,\lambda \tau} = 0$  under the null hypothesis of  $\lambda = 0$ , we can apply the central limit theorem in Lemma 3 for regression (17) to obtain  $\sqrt{T} R^{Ind,\lambda \hat{\tau}} = \sqrt{T} R^{Ind,\lambda}(\hat{\tau} - \tau) \rightarrow_d N(0, R^{Ind,\lambda} \Sigma_\tau R^{Ind,\lambda'})$ . Furthermore, by the consistency of  $\hat{\Sigma}_\tau$  we have  $F_T^{Ind,\lambda} \rightarrow_d \chi_1^2$ . Note that if  $e_{t+1}$  is homoskedastic, then the proof is analogous to the heteroskedastic case, except  $\hat{\Sigma}_\tau \rightarrow_p \sigma_e^2 \Sigma_{zz}$  which can be consistently estimated by  $\hat{\sigma}_e^2 \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_t \hat{z}_t' \right)^{-1}$ , with  $\hat{\sigma}_e^2 = (1/T) \sum_{t=1}^{T-1} \hat{e}_{t+1}^2$  a consistent estimator of  $\sigma_e^2$  and  $\hat{e}_{t+1}$  the least squares residuals. ■

**Proof of Theorem 2:** The proof of  $F_T^{Ind,\varsigma} \rightarrow_d \chi_h^2$  for testing  $\varsigma_1 = \dots = \varsigma_{\hat{q}} = 0$  in regression (20) follows immediately from Lemma 3 and the proof of Theorem 1. We now establish the asymptotic chi-squared distribution for the test of the null hypothesis  $\dot{H}_0 : \psi_1 = \dots = \psi_{\hat{q}} = 0$  in regression (19). For simplicity of exposition, let  $\dot{p} = \dot{q} = \dot{h} = 1$ . Following the steps in the proof of Theorem 1, we note that the infeasible regression (19) can be rewritten as

$$H^{-1} \hat{f}_{t+1} = \hat{z}_t' \tau + u_{t+1} + \rho H^{-1} (H f_t - \hat{f}_t) - H^{-1} (H f_{t+1} - \hat{f}_{t+1}).$$

In the following, we denote  $\tau = (\nu, \kappa, \psi, \rho H^{-1})'$  and  $\hat{z}_t = (1, X_t, Y_t, \hat{f}_t)'$ . It is important to remark that, comparing with the standard set up in Theorem 1, the above expression has an extra term  $H^{-1} (H f_{t+1} - \hat{f}_{t+1})$ , because the dependent variable also has to be replaced by the estimated factor  $\hat{f}_{t+1}$  in the feasible regression. We now write the model in matrix form

$$\hat{f}^1 H^{-1'} = \hat{z} \tau + u + (f H' - \hat{f}) H^{-1'} \rho - (f^1 H' - \hat{f}^1) H^{-1'}, \quad (48)$$

where  $\hat{f}^1 = (\hat{f}_2, \dots, \hat{f}_T)'$ ,  $f^1 = (f_2, \dots, f_T)'$ ,  $\hat{f} = (\hat{f}_1, \dots, \hat{f}_{T-1})'$  and  $f = (f_1, \dots, f_{T-1})'$ ,  $u = (u_2, \dots, u_T)'$  and  $\hat{z} = (\hat{z}_1, \dots, \hat{z}_{T-1})'$ . Hence, least squares estimation of (48) yields

$$\begin{aligned} \sqrt{T}(\hat{\tau} - \tau) &= (T^{-1}\hat{z}'\hat{z})^{-1} T^{-1/2}\hat{z}'u + (T^{-1}\hat{z}'\hat{z})^{-1} \left[ T^{-1/2}\hat{z}' \left( fH' - \hat{f} \right) \right] H^{-1'} \rho \\ &\quad - (T^{-1}\hat{z}'\hat{z})^{-1} \left[ T^{-1/2}\hat{z}' \left( f^1H' - \hat{f}^1 \right) \right] H^{-1'}. \end{aligned}$$

The second term on the right-hand side of last equation is  $O_p(T^{1/2}/\min(N, T)) = o_p(1)$  by the result (iii) of Lemma 1 if  $\sqrt{T}/N \rightarrow 0$ . In addition, the third term is  $T^{-1/2} \sum_{t=1}^T (\hat{f}_{t+1} - Hf_{t+1})\hat{z}'_t = O_p(T^{1/2}/\min(N, T)) = o_p(1)$  when  $\sqrt{T}/N \rightarrow 0$  according to (ii) of Lemma 2.

By the result (iv) of Lemma 1 and  $\sqrt{T}/N \rightarrow 0$ ,  $T^{-1/2}\hat{z}'^{-1/2}\Phi z'u + o_p(1)$ , with  $\Phi = \text{diag}(I, H)$  a block diagonal matrix, we obtain that  $\sqrt{T}(\hat{\tau} - \tau) = (T^{-1}\hat{z}'\hat{z})^{-1}\Phi T^{-1/2}z'u + o_p(1)$  and  $\sqrt{T}(\hat{\tau} - \tau) \rightarrow_d N(0, \Sigma_\tau)$ , where  $\Sigma_\tau$  can be consistently estimated using

$$\hat{\Sigma}_\tau = \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_t \hat{z}'_t \right)^{-1} \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{u}_{t+1}^2 \hat{z}_t \hat{z}'_t \right) \left( \frac{1}{T} \sum_{t=1}^{T-1} \hat{z}_t \hat{z}'_t \right)^{-1}$$

with  $\hat{u}_{t+1}$  the OLS residuals. Finally, let  $R^{Ind, \psi} = (0, 0, 1, 0)$ . The rest of the proof for the  $F$ -statistic  $F_T^{Ind, \psi} = \left( \sqrt{T} R^{Ind, \psi} \hat{\tau} \right) \left( R^{Ind, \psi} \hat{\Sigma}_\tau R^{Ind, \psi'} \right)^{-1} \left( \sqrt{T} R^{Ind, \psi} \hat{\tau} \right)' \rightarrow_d \chi_1^2$  follows straightforwardly from the results on the asymptotic normality of  $\hat{\tau}$ ,  $\sqrt{T} R^{Ind, \psi} \hat{\tau} = \sqrt{T} R^{Ind, \psi} (\hat{\tau} - \tau) \rightarrow_d N(0, R^{Ind, \psi} \Sigma_\tau R^{Ind, \psi'})$  under the null hypothesis of  $\psi = 0$ , and the consistency of  $\hat{\Sigma}_\tau$  to  $\Sigma_\tau$ . ■

## B Appendix: Data

Table 17: Big data: Description of the variables

id	tcode	fred	description
<b>Group 1: Output and income</b>			
1	1	5	RPI Real Personal Income
2	2	5	W875RX1 Real personal income ex transfer receipts
3	6	5	INDPRO IP Index
4	7	5	IPFPNSS IP: Final Products and Nonindustrial Supplies
5	8	5	IPFINAL IP: Final Products (Market Group)
6	9	5	IPCONGD IP: Consumer Goods
7	10	5	IPDCONGD IP: Durable Consumer Goods
8	11	5	IPNCONGD IP: Nondurable Consumer Goods
9	12	5	IPBUSEQ IP: Business Equipment
10	13	5	IPMAT IP: Materials
11	14	5	IPDMAT IP: Durable Materials
12	15	5	IPNMAT IP: Nondurable Materials
13	16	5	IPMANSICS IP: Manufacturing (SIC)
14	17	5	IPB51222s IP: Residential Utilities
15	18	5	IPFUELS IP: Fuels
16	19	1	NAPMPI ISM Manufacturing: Production Index
17	20	2	CUMFNS Capacity Utilization: Manufacturing
<b>Group 2: Labor market</b>			
1	21	2	HWI Help-Wanted Index for United States
2	22	2	HWIURATIO Ratio of Help Wanted/No. Unemployed
3	23	5	CLF16OV Civilian Labor Force
4	24	5	CE16OV Civilian Employment
5	25	2	UNRATE Civilian Unemployment Rate
6	26	2	UEMPMEAN Average Duration of Unemployment (Weeks)
7	27	5	UEMPLT5 Civilians Unemployed - Less Than 5 Weeks
8	28	5	UEMP5TO14 Civilians Unemployed for 5-14 Weeks
9	29	5	UEMP15OV Civilians Unemployed - 15 Weeks & Over
10	30	5	UEMP15T26 Civilians Unemployed for 15-26 Weeks

**Note:** This table presents the variables included in the groups “Output and income” and “Labor market”.

Table 18: Big data: Description of the variables (Cont.)

id	tcode	fred		description
<b>Group 2: Labor market (Cont.)</b>				
11	31	5	UEMP27OV	Civilians Unemployed for 27 Weeks and Over
12	32	5	CLAIMSx	Initial Claims
13	33	5	PAYEMS	All Employees: Total nonfarm
14	34	5	USGOOD	All Employees: Goods-Producing Industries
15	35	5	CES1021000001	All Employees: Mining and Logging: Mining
16	36	5	USCONS	All Employees: Construction
17	37	5	MANEMP	All Employees: Manufacturing
18	38	5	DMANEMP	All Employees: Durable goods
19	39	5	NDMANEMP	All Employees: Nondurable goods
20	40	5	SRVPRD	All Employees: Service-Providing Industries
21	41	5	USTPU	All Employees: Trade, Transportation & Utilities
22	42	5	USWTRADE	All Employees: Wholesale Trade
23	43	5	USTRADE	All Employees: Retail Trade
24	44	5	USFIRE	All Employees: Financial Activities
25	45	5	USGOVT	All Employees: Government
26	46	1	CES0600000007	Avg Weekly Hours : Goods-Producing
27	47	2	AWOTMAN	Avg Weekly Overtime Hours : Manufacturing
28	48	1	AWHMAN	Avg Weekly Hours : Manufacturing
29	49	1	NAPMEI	ISM Manufacturing: Employment Index
30	127	6	CES0600000008	Avg Hourly Earnings : Goods-Producing
31	128	6	CES2000000008	Avg Hourly Earnings : Construction
32	129	6	CES3000000008	Avg Hourly Earnings : Manufacturing
<b>Group 3: Housing</b>				
1	50	4	HOUST	Housing Starts: Total New Privately Owned
2	51	4	HOUSTNE	Housing Starts, Northeast
3	52	4	HOUSTMW	Housing Starts, Midwest
4	53	4	HOUSTS	Housing Starts, South
5	54	4	HOUSTW	Housing Starts, West

**Note:** This table presents the variables included in the groups “Labor market” and “Housing”.

Table 19: Big data: Description of the variables (Cont.)

id	tcode	fred	description
<b>Group 3: Housing (Cont.)</b>			
6	55	4	PERMIT New Private Housing Permits (SAAR)
7	56	4	PERMITNE New Private Housing Permits, Northeast (SAAR)
8	57	4	PERMITMW New Private Housing Permits, Midwest (SAAR)
9	58	4	PERMITS New Private Housing Permits, South (SAAR)
10	59	4	PERMITW New Private Housing Permits, West (SAAR)
<b>Group 4: Consumption, orders, and inventories</b>			
1	3	5	DPCERA3M086SBEA Real personal consumption expenditures
2	4	5	CMRMTSPLx Real Manu. and Trade Industries Sales
3	5	5	RETAILx Retail and Food Services Sales
4	60	1	NAPM ISM : PMI Composite Index
5	61	1	NAPMNOI ISM : New Orders Index
6	62	1	NAPMSDI ISM : Supplier Deliveries Index
7	63	1	NAPMII ISM : Inventories Index
8	64	5	ACOGNO New Orders for Consumer Goods
9	65	5	AMDMNOx New Orders for Durable Goods
10	66	5	ANDENOx New Orders for Nondefense Capital Goods
11	67	5	AMDMUOx Unfilled Orders for Durable Goods
12	68	5	BUSINVx Total Business Inventories
13	69	2	ISRATIOx Total Business: Inventories to Sales Ratio
14	130	2	UMCSENTx Consumer Sentiment Index
<b>Group 5: Money and credit</b>			
1	70	6	M1SL M1 Money Stock
2	71	6	M2SL M2 Money Stock
3	72	5	M2REAL Real M2 Money Stock
4	73	6	AMBSL St. Louis Adjusted Monetary Base
5	74	6	TOTRESNS Total Reserves of Depository Institutions
6	75	7	NONBORRES Reserves Of Depository Institutions
7	76	6	BUSLOANS Commercial and Industrial Loans

**Note:** This table presents the variables included in the groups “Housing”, “Consumption, orders, and inventories” and “Money and credit”.

Table 20: Big data: Description of the variables (Cont.)

id	tcode	fred	description
<b>Group 5: Money and credit (Cont.)</b>			
8	77	6	REALLN Real Estate Loans at All Commercial Banks
9	78	6	NONREVSL Total Nonrevolving Credit
10	79	2	CONSPI Nonrevolving consumer credit to Personal Income
11	131	6	MZMSL MZM Money Stock
12	132	6	DTCOLNVHFNM Consumer Motor Vehicle Loans Outstanding
13	133	6	DTCTHFNM Total Consumer Loans and Leases Outstanding
14	134	6	INVEST Securities in Bank Credit at All Commercial Banks
<b>Group 6: Interest and exchange rates</b>			
1	84	2	FEDFUNDS Effective Federal Funds Rate
2	85	2	CP3Mx 3-Month AA Financial Commercial Paper Rate
3	86	2	TB3MS 3-Month Treasury Bill:
4	87	2	TB6MS 6-Month Treasury Bill:
5	88	2	GS1 1-Year Treasury Rate
6	89	2	GS5 5-Year Treasury Rate
7	90	2	GS10 10-Year Treasury Rate
8	91	2	AAA Moody's Seasoned Aaa Corporate Bond Yield
9	92	2	BAA Moody's Seasoned Baa Corporate Bond Yield
10	93	1	COMPAPFFx 3-Month Commercial Paper Minus FEDFUNDS
11	94	1	TB3SMFFM 3-Month Treasury C Minus FEDFUNDS
12	95	1	TB6SMFFM 6-Month Treasury C Minus FEDFUNDS
13	96	1	T1YFFM 1-Year Treasury C Minus FEDFUNDS
14	97	1	T5YFFM 5-Year Treasury C Minus FEDFUNDS
15	98	1	T10YFFM 10-Year Treasury C Minus FEDFUNDS
16	99	1	AAAFFM Moody's Aaa Corporate Bond Minus FEDFUNDS
17	100	1	BAAFFM Moody's Baa Corporate Bond Minus FEDFUNDS
18	101	5	TWEXMMTH Trade Weighted U.S. Dollar Index: Major Currencies
19	102	5	EXSZUSx Switzerland / U.S. Foreign Exchange Rate
20	103	5	EXJPUSx Japan / U.S. Foreign Exchange Rate

**Note:** This table presents the variables included in the groups “Money and credit” and “Interest and exchange rates”.



Table 21: Big data: Description of the variables (Cont.)

id	tcode	fred	description
<b>Group 6: Interest and exchange rates (Cont.)</b>			
21	104	5	EXUSUKx U.S. / U.K. Foreign Exchange Rate
22	105	5	EXCAUSx Canada / U.S. Foreign Exchange Rate
<b>Group 7: Prices</b>			
1	106	6	WPSFD49207 PPI: Finished Goods
2	107	6	WPSFD49502 PPI: Finished Consumer Goods
3	108	6	WPSID61 PPI: Intermediate Materials
4	109	6	WPSID62 PPI: Crude Materials
5	110	6	OILPRICEx Crude Oil, spliced WTI and Cushing
6	111	6	PPICMM PPI: Metals and metal products:
7	112	1	NAPMPRI ISM Manufacturing: Prices Index
8	113	6	CPIAUCSL CPI : All Items
9	114	6	CPIAPPSL CPI : Apparel
10	115	6	CPITRNSL CPI : Transportation
11	116	6	CPIMEDSL CPI : Medical Care
12	117	6	CUSR0000SAC CPI : Commodities
13	118	6	CUUR0000SAD CPI : Durables
14	119	6	CUSR0000SAS CPI : Services
15	120	6	CPIULFSL CPI : All Items Less Food
16	121	6	CUUR0000SA0L2 CPI : All items less shelter
17	122	6	CUSR0000SA0L5 CPI : All items less medical care
18	123	6	PCEPI Personal Cons. Expend.: Chain Index
19	124	6	DDURRG3M086SBEA Personal Cons. Exp: Durable goods
20	125	6	DNDGRG3M086SBEA Personal Cons. Exp: Nondurable goods

**Note:** This table presents the variables included in the groups “Interest and exchange rates” and “Prices”.

Table 22: Big data: Description of the variables (Cont.)

id	tcode	fred	description
<b>Group 7: Prices (Cont.)</b>			
21	126	6	DSERRG3M086SBEA Personal Cons. Exp: Services
<b>Group 8: Stock market</b>			
1	80	5	S&P 500 S&P’s Common Stock Price Index: Composite
2	81	5	S&P: indust S&P’s Common Stock Price Index: Industrials
3	82	2	S&P div yield S&P’s Composite Common Stock: Dividend Yield
4	83	5	S&P PE ratio S&P’s Composite Common Stock: Price-Earnings Ratio
5	135	1	VXOCLSx VXO

**Note:** This table presents the variables included in the groups “Prices” and “Stock market”.