

1 RESEARCH

2 **Non-parametric test for connectivity detection in multivariate**
3 **autoregressive networks and application to multiunit activity data**

4 M Gilson^{1,5}, A Tauste Campo^{1,2,5}, X Chen³, A Thiele³, and G Deco^{1,4}

5 ¹Computational Neuroscience Group, Dept. de Tecnologies de la Informació i les Comunicacions, Universitat Pompeu Fabra, Barcelona, Spain

6 ²Epilepsy Monitoring Unit, Department of Neurology, Hospital del Mar Medical Research Institute, Barcelona, Spain

7 ³Institute of Neuroscience, Newcastle University, Newcastle upon Tyne, UK

8 ⁴Institució Catalana de Recerca i Estudis Avançats, Barcelona, Spain

9 ⁵equal contribution

10 **Keywords:** Network connectivity detection; non-parametric significance method; multivariate autoregressive process;
11 Granger causality; multiunit activity

ABSTRACT

12 Directed connectivity inference has become a cornerstone in neuroscience to analyze multivariate data
13 from neuroimaging and electrophysiological techniques. Here we propose a non-parametric significance
14 method to test the non-zero values of multivariate autoregressive model to infer interactions in recurrent
15 networks. We use random permutations or circular shifts of the original time series to generate the
16 null-hypothesis distributions. The underlying network model is the same as used in multivariate Granger
17 causality, but our test relies on the autoregressive coefficients instead of error residuals. By means of
18 numerical simulation over multiple network configurations, we show that this method achieves a good
19 control of false positives - type 1 error - and detects existing pairwise connections more accurately than
20 using the standard parametric test for the ratio of error residuals. In practice, our method aims to detect
21 temporal interactions in real neuronal networks with nodes possibly exhibiting redundant activity. As a
22 proof of concept, we apply our method to multiunit activity (MUA) recorded from Utah electrode arrays
23 in a monkey and examine detected interactions between 25 channels. We show that during stimulus
24 presentation our method detects a large number of interactions that cannot be solely explained by the
25 increase in the MUA level.

INTRODUCTION

26 In recent years, there has been a growing interest in developing multivariate techniques to infer causal
27 relations among time series. The initial formulation of the problem goes back to the seminal work by
28 Granger in 1960's (Granger, 1969) motivated by the analysis of the pairwise influence between economic
29 time series. In this work, Granger decomposes the cross-spectrum of two autoregressive time series into
30 two directional components that account for the potential causal influences between each other. A general
31 solution of the problem in multivariate scenarios was developed a decade later by the introduction of
32 multivariate autoregressive (MVAR) processes, which allow the estimation of causal relationships
33 between nodes in networks with linear feedback based on their observed activity (Amemiya, 1974;
34 Geweke, 1982, 1984; Lütkepohl, 2005). The MVAR was further combined with spectral analysis to
35 develop the directed transfer entropy function (Kaminski & Blinowska, 1991; Kamiński, Ding, Truccolo,
36 & Bressler, 2001), which has been employed to analyze connectivity patterns in neurobiological systems
37 (Babiloni et al., 2005; Wilke, Ding, & He, 2008). Granger causality analysis is nowadays often used to
38 evaluate the influence of a group of variables onto other, which corresponds to the influence of a subgroup
39 of nodes onto another one in networks. Nevertheless, it is also applied to detect individual connections
40 between pairs of nodes (each subgroup being a single node), which sets the context of the present paper.

41 In neuroscience, this inference problem has been transposed to analyze interactions between neuronal
42 populations from spiking activity or neuroimaging measurements such as fMRI, EEG and MEG (Lusch,
43 Maia, & Kutz, 2016; Messé, Rudrauf, Benali, & Marrelec, 2014; Michalareas, Schoffelen, Paterson, &
44 Gross, 2013; Rogers, Katwal, Morgan, Asplund, & Gore, 2010; Seth, Barrett, & Barnett, 2015; Storkey et
45 al., 2007). Two types of estimation procedures may be distinguished: measures relying on an underlying
46 interaction model such as Granger causality analysis (M. Ding, Chen, & Bressler, 2006) and dynamic
47 causal modeling (DCM) (Friston, Harrison, & Penny, 2003) on the one hand; and model-free measures
48 such as transfer entropy (Schreiber, 2000) and directed information (Massey, 1990), which make
49 minimal model assumptions on the other hand. Although model-free approaches have proven useful to
50 describe neural propagation at spike-train level (So, Koralek, Ganguly, Gastpar, & Carmena, 2012;
51 Tauste Campo et al., 2015), certain assumptions are required when estimating interactions at the neuronal

52 population level, in which broader spatial and temporal scales contribute to shaping the signals.
53 Motivated by data-driven practical problems, methodological refinements of Granger causality analysis
54 (or MVAR-based methods) have considered additive noise (Vinck et al., 2015) or measurement noise via
55 state-space models (Barnett & Seth, 2015; Friston et al., 2014). However, in the majority of cases, the
56 ratio behind the detection test concerns sub-model error residuals, which might become too similar when
57 connections are placed in a highly redundant network, thus increasing the missed detection rate
58 (Stramaglia, Cortes, & Marinazzo, 2014).

59 To overcome the difficulties of detecting directed connections in the general context of large networks,
60 we propose to test the significance of the MVAR coefficients using a non-parametric procedure. As a
61 generative model, the MVAR process is canonically related to Granger causality analysis: the linear
62 regression in the upper right inset of Fig. 1A provides both coefficients and residuals, the latter being
63 viewed as the remaining uncertainty in the prediction of the target time series by its source(s). By
64 comparing the residuals of two linear regressions - one involving a supposed driver node and one without
65 it - in a log ratio, traditional tests for Granger causality estimate the effective interaction of one node onto
66 another (Barrett & Barnett, 2013). Since these log ratios asymptotically converge to known distributions,
67 parametric statistical tests have been developed to assess the significance of these interactions (Barnett &
68 Seth, 2014). Instead, our proposed method evaluates the significance of the MVAR coefficients to infer
69 the existence of network connections. To achieve this, we propose a non-parametric significance test in
70 the regression coefficients space. Previous literature on non-parametric testing for Granger causality has
71 resorted to surrogate data generated by trial shuffling (Dhamala, Rangarajan, & Ding, 2008; Nedungadi,
72 Rangarajan, Jain, & Ding, 2009), bootstrap procedures (Diks & DeGoede, 2001) or by phase
73 randomization in frequency-domain measures (L. Ding, Worrell, Lagerlund, & He, 2007; Li et al., 2016).
74 Here we focus on within-trial surrogate tests for time-domain coefficients and compare them across
75 standard techniques (Faes, Marinazzo, Montalto, & Nollo, 2014; Schreiber & Schmitz, 1996; Winkler,
76 Ridgway, Webster, Smith, & Nichols, 2014).

77 Our approach is motivated by the growing of multichannel recording techniques in neuroscience,
78 which require tailored multivariate analysis. In the context of recurrent networks, which are ubiquitous in
79 neuroscience, we provide numerical evidence that these tests can achieve a good control of the
80 false-alarm rate and might improve the miss rate by properly adapting the null distribution to each

81 connection. The focus of the present analysis is on the case where we observe more time samples (a few
 82 thousands per node) than the network size (about a hundred nodes), a usual ground for
 83 electrophysiological data. Within this regime, we test the robustness of the detection method for a broad
 84 range of network parameters and various non-trivial topologies inspired by neuronal networks.

METHODS: MULTIVARIATE AUTOREGRESSIVE MODEL AND CONNECTIVITY ESTIMATION

The activity in the MVAR process - a.k.a. noise-diffusion discrete-time network - is described by the following equation:

$$x^t = Ax^{t-1} + \zeta^t, \quad (1)$$

85 where the connection matrix A describes the interactions between coordinates of the vector $x^t = (x_i^t)$
 86 with time t being an integer and node index $1 \leq i \leq N$. Here we constrain our study to the case where ζ^t
 87 is Gaussian (possibly cross-correlated noise), whose realizations are time independent for successive time
 88 steps. Without loss of generality, we assume that all variables ζ^t have zero means, giving the same
 89 property for all x_i . We only consider MVAR processes of order 1 in a first place, but will extend the work
 90 to the case of order 2 in a later section.

91 *Granger causality analysis*

Granger causality analysis is usually presented using time series and the estimation of non-zero coefficients in A from observed activity over a period $1 \leq t \leq T$ relies on the linear regression of the activity x_i^t of a given node i at time t by the past activity of a subset \mathcal{S} of network nodes:

$$x_i^t = \sum_{j \in \mathcal{S}} A_{ij} x_j^{t-1} + \epsilon^t \quad (2)$$

for $2 \leq t \leq T$. When T is large, the coefficients a_{ij} converge toward A_{ij} . We define the residual ϵ as the standard deviation of the ϵ^t for the ordinary least-square (OLS) regression in Eq. (2), which is for

$$\epsilon(x_i^{2 \leq t \leq T} | x_{\mathcal{S}}^{1 \leq t \leq T-1}) = \sqrt{\sum_t (\epsilon^t)^2}, \quad (3)$$

92 with a notation similar to conditional probabilities; the superscript t indicates the considered time range
 93 and the subscripts indicate the nodes involved. To detect the existence of connection $j \rightarrow i$ in a network,
 94 two types of Granger causality analysis exist: ‘unconditional’ and ‘conditional’ (Geweke, 1982, 1984). ,

95 they consider the comparisons of the following residuals:

$$\begin{aligned} \text{GRu}(x_j \rightarrow x_i) &= \ln \left[\frac{\epsilon(x_i^{2 \leq t \leq T} | x_i^{1 \leq t \leq T-1})}{\epsilon(x_i^{2 \leq t \leq T} | x_{i,j}^{1 \leq t \leq T-1})} \right]; \\ \text{GRc}(x_j \rightarrow x_i) &= \ln \left[\frac{\epsilon(x_i^{2 \leq t \leq T} | x_{1, \dots, j-1, j+1, \dots, N}^{1 \leq t \leq T-1})}{\epsilon(x_i^{2 \leq t \leq T} | x_{1, \dots, N}^{1 \leq t \leq T-1})} \right]. \end{aligned} \quad (4)$$

For both GRu and GRc, which have a univariate target node x_i , the usual parametric test for significance relies on the F statistics, which performs better for small number of samples (Barnett & Seth, 2014). The null hypothesis of no interaction for $\text{GRu}(x_j \rightarrow x_i)$ corresponds to $m = T$, $p = 1$, $n_x = 1$ and $n_y = 2$ using the notation in Barnett and Seth (2014)

$$\frac{\epsilon(x_i^{2 \leq t \leq T} | x_i^{1 \leq t \leq T-1}) - \epsilon(x_i^{2 \leq t \leq T} | x_{i,j}^{1 \leq t \leq T-1})}{\epsilon(x_i^{2 \leq t \leq T} | x_{i,j}^{1 \leq t \leq T-1})} = [\exp(\text{GRu}_{ij}) - 1] > \frac{\phi(\alpha, 1, T-3)}{T-3} \quad (5)$$

with α the desired sensitivity and ϕ the inverse survival function of the F-distribution (www.scipy.org, n.d.). The equivalent for GRc corresponds to $n_y = N$, yielding

$$\frac{\epsilon(x_i^{2 \leq t \leq T} | x_{1, \dots, N}^{1 \leq t \leq T-1}) - \epsilon(x_i^{2 \leq t \leq T} | x_{1, \dots, j-1, j+1, \dots, N}^{1 \leq t \leq T-1})}{\epsilon(x_i^{2 \leq t \leq T} | x_{1, \dots, j-1, j+1, \dots, N}^{1 \leq t \leq T-1})} > \frac{\phi(\alpha, 1, T-N-1)}{T-N-1}. \quad (6)$$

96 We also use non-parametric tests for GRc by performing a circular shift (see details below in Section
97 ‘Generation of surrogate time series’) either on the target node for each connection (Faes et al., 2014) or
98 independently on the time series of all nodes (in order to save time in estimating the full network’s
99 connectivity by shuffling somehow all targets simultaneously). Both cases provide a null distribution for
100 the log ratio, with which the actual estimated log ratio can be compared.

101 *Multivariate autoregressive (MVAR) estimation*

To detect the existence of connections $A_{ij} > 0$, another possibility is to estimate the coefficients themselves, which can be done using the covariances of the observed activity variables x^t (Lütkepohl, 2005):

$$\widehat{Q}_{ij}^\tau = \frac{1}{T - \tau_{\max} - 1} \sum_{1 \leq t \leq T - \tau_{\max}} (x_i^{t+\tau} - \bar{x}_i)(x_j^t - \bar{x}_j), \quad (7)$$

where T denotes the number of successive samples indexed by t , $\tau \in \{0, 1\}$ is the time shift (here $\tau_{\max} = 1$) and the observed mean activity for each node is $\bar{x}_i = \frac{1}{T} \sum_t x_i^t$. The Yule-Walker equation gives a consistency equation for the theoretical covariance matrices (without hat) in terms of the connectivity A

in the dynamics described by Eq.(1):

$$Q^1 = A Q^0 . \quad (8)$$

The estimation of network connections relies on evaluating A from Eq. (8) for the empirical covariance matrices defined as Eq. (7) and calculated for a given time series:

$$A = \widehat{Q}^1 (\widehat{Q}^0)^{-1} . \quad (9)$$

102 Note that this OLS estimate corresponds to the linear regression related to $\epsilon(x_{1,\dots,N}^{2 \leq t \leq T} | x_{1,\dots,N}^{1 \leq t \leq T-1})$ and is also
103 to the linear model with maximum likelihood under the assumption that the observed process is Gaussian.

104 *MVAR of order 2*

Eq. (1) can be extended to the case where the activity vector x^t is determined by the two previous time steps:

$$x^t = A^1 x^{t-1} + A^2 x^{t-2} + \zeta^t . \quad (10)$$

For the second order, we use $\tau_{\max} = 2$ in Eq. (7) and the estimation of A^1 and A^2 via the Yule -Walker equation is given by (Lütkepohl, 2005, p. 86)

$$\tilde{A} = \tilde{Q}^1 (\tilde{Q}^0)^{-1}, \quad (11)$$

105 with the block matrices

$$\begin{aligned} \tilde{A} &= \begin{pmatrix} A^1 & A^2 \end{pmatrix}, \\ \tilde{Q}^0 &= \begin{pmatrix} \widehat{Q}^0 & \widehat{Q}^1 \\ (\widehat{Q}^1)^\dagger & \widehat{Q}^0 \end{pmatrix}, \\ \tilde{Q}^1 &= \begin{pmatrix} \widehat{Q}^1 & \widehat{Q}^2 \end{pmatrix}. \end{aligned} \quad (12)$$

106 The coefficients of A^1 and A^2 can thus be estimated using a matrix multiplication and an inversion
107 involving the covariances, as with the first-order case in Eq. (9).

108 *Generation of surrogate time series*

109 In this paper, we consider circular shifts (CS), random permutations (RP) and phase randomization (PR)
110 to shuffle the time points of the observed time series. From the original x_i^t with $1 \leq t \leq T$,

- 111 ▪ CS draws a random integer $t_0 \in \{1, \dots, T\}$ and returns $(x_i^{t_0}, \dots, x_i^T, x_i^1, \dots, x_i^{t_0-1})$;
- 112 ▪ RP draws a random permutation σ of $\{1, \dots, T\}$ such that each integer appears once (and only
113 once) and returns $x_i^{\sigma(t)}$;
- 114 ▪ PR calculates the discrete Fourier transform $\mathcal{F}(x_i^t)$ of the original x_i^t , then multiplies each of the T
115 coefficients of $\mathcal{F}(x_i^t)$ by $\exp(2\pi i \phi^t)$ with ϕ^t randomly chosen in $[0, 2\pi]$, and performs the inverse
116 transform.

117 Importantly, these operations are applied to each time series independently of the others.

118 In addition, we consider the replacement of all time series in the network by T normally distributed
119 variables with a standard deviation equal to the mean of the standard deviations of x_i^t along the time axis,
120 then averaged for all nodes. We refer to these surrogates as STD.

121 *Experimental setup and processing of electrode measurements to extract MUAe activity*

122 All procedures were carried out in accordance with the European Communities Council Directive RL
123 2010/63/EC, the US National Institutes of Health Guidelines for the Care and Use of Animals for
124 Experimental Procedures, and the UK Animals Scientific Procedures Act. Two male macaque monkeys
125 (5 - 14 years of age) were used in the experiment; only the data for the first one is used here. A surgical
126 operation was performed under sterile conditions, in which a custom-made head post (Peek, Tecapeek)
127 was embedded into a dental acrylic head stage. Details of surgical procedures and post-operative care
128 have been published previously (Thiele, Delicato, Roberts, & Gieselmann, 2006). During the surgery
129 microelectrode chronic Utah arrays (5*5 grids), attached to a CerePort base (Blackrock Microsystems)
130 were implanted into V1. Electrodes were 1 mm in length in line with procedures described in Supèr and
131 Roelfsema (2005). Stimulus presentation was controlled using CORTEX software (Laboratory of
132 Neuropsychology, NIMH, <http://dally.nimh.nih.gov/index.html>) on a computer with an Intel Core i3-540
133 processor. Stimuli were displayed at a viewing distance of 0.54 m, on a 25" Sony Trinitron CRT monitor
134 with a resolution of 1280 by 1024 pixels, yielding a resolution of 31.5 pixels / degree of visual angle
135 (dva). The monitor refresh rate was 85 Hz for monkey 1, and 75 Hz for monkey 2. A gamma correction
136 was used to linearize the monitor output, and the gratings had 50% contrast. Monkeys performed a
137 passive viewing task where they fixated centrally while stationary sinusoidal grating of either horizontal
138 or vertical orientation and 2 cycle per degree spatial frequency, were presented in a location that covered

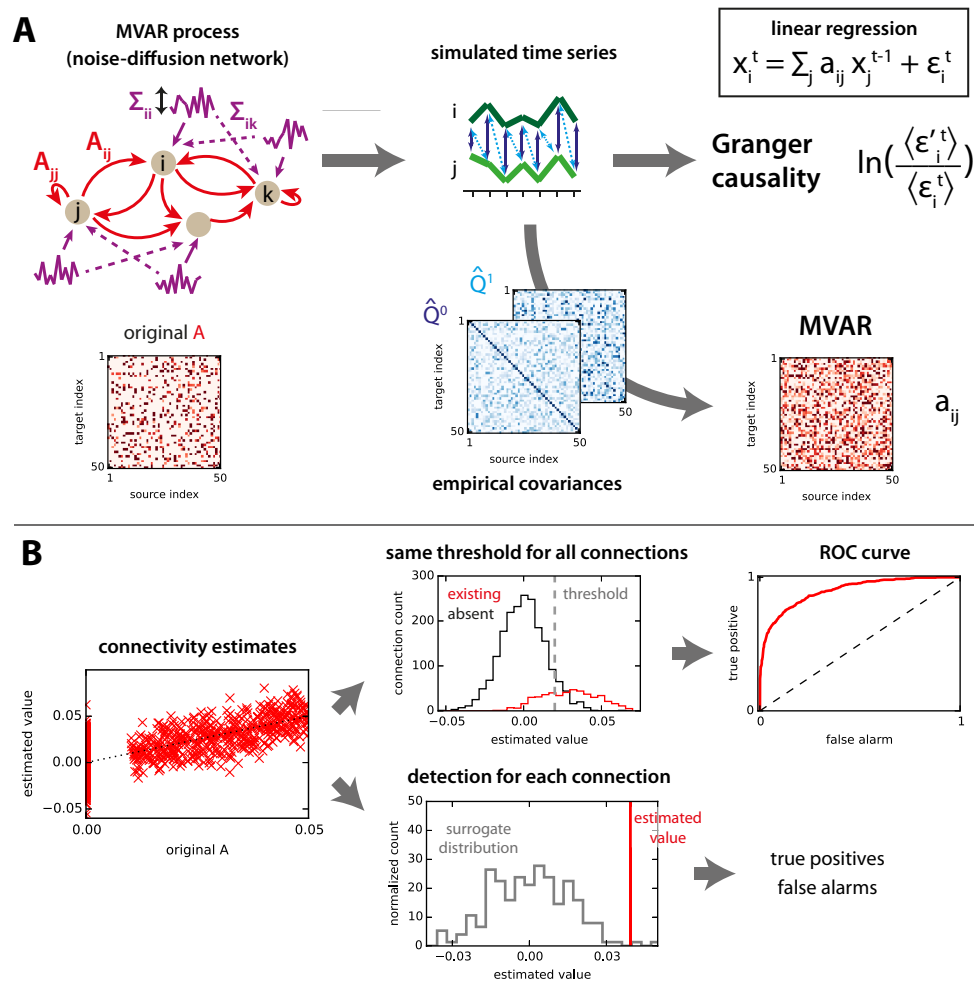
139 all receptive fields recorded from the 25 electrode tips. Stimuli were presented 500 ms after fixation onset
140 for 150 ms. Raw data were acquired at a sampling frequency of 32556 Hz using a 64-channel Digital
141 Lynx 16SX Data Acquisition System (Neuralynx, Inc.). Following each recording session, the raw data
142 were processed offline using commercial (Neuralynx, Inc.). Signals were extracted using Cheetah 5 Data
143 Acquisition Software, with bandpass filtering set to allow for spike extraction (600-9000 Hz) and saved at
144 16-bit resolution.

145 In the present study, we focus on the period starting 200 ms before and finishing 200 ms after the
146 stimulus onset, for 4 conditions (vertical gradings with pre/post cue in the receptor/opposite field) that
147 will not be compared in details. The electrode recordings is firstly down-sampled from 32556 Hz to
148 1000 Hz. A high-pass filter above 400 Hz is then applied - 3rd-order Butterworth filter at 0.8 of the
149 Nyquist frequency (www.scipy.org, n.d.) - followed by a smoothing of 4 ms to extract the envelope of the
150 resulting signal, by retaining the 250 time points of 1000-ms period surrounding the stimulus onset.

BENCHMARK OF DETECTION PERFORMANCE FOR SYNTHETIC DATA

159 The workflow of the benchmark for the estimation procedure is schematically represented in Fig. 1A. We
160 first consider a MVAR process defined by Eq. (1) with given connectivity matrix A and input covariances
161 (obtained by mixing independent Gaussian processes) to generate the activity of the network. From the
162 observed activity over a period of duration T , we estimate the coefficients matrix A using the covariances
163 as described in Eq. (9). We also perform the linear regressions of each node activity over the past activity
164 of given subsets of nodes corresponding to the unconditional (GRu) and conditional (GRc) Granger
165 causality analysis, from which we calculate the ratios of residuals in Eq. (4). Actually, these estimates
166 correspond to the same ordinary least-square (OLS) regression (top right in Fig. 1A) and the difference
167 resides in the spaces where they lie: coefficients versus residuals.

168 For each method, the prediction power can be measured by the relationship between the estimated
169 values and the original connectivity values, as illustrated in Fig. 1B (left). To discriminate between
170 existing and absent connections, one can apply a common threshold for all connections (top thread); by
171 sliding this threshold, we obtain the ROC curve with the rate of false alarms on the x-axis and true
172 positives on the y-axis. The area under the curve indicates in a single value how well the ranking of
173 estimated value performs for the detection of connections in the original connectivity. Alternatively, an



151 **Figure 1. Network model and connectivity estimation.** **A:** For a given directed connectivity A and input covariances Σ (left), the network activity
 152 (middle) is simulated using Eq. (1). From the observed time series, the existing interactions in the original connectivity can be estimated (right): Granger
 153 causality analysis uses the residuals of linear regressions (ϵ in the upper right equation; see Methods for details about the residuals used in the log ratio),
 154 whereas MVAR corresponds to the coefficients. Note that MVAR can be obtained using the empirical covariance matrices \hat{Q}^0 and \hat{Q}^1 , see in Eq. (7) for $\tau = 0$
 155 and 1. **B:** The left panel compares the estimated values to the original values for all connections in the network. The upper thread displays the distributions
 156 of estimated values for existing and non-existing connections in the original network. Using a sliding threshold (vertical dashed gray line) on the estimated
 157 values, one can calculate the ROC curve (right). The lower thread compares the estimated value for a single connection to a null distribution. From this, the
 158 choice whether the connection exists or not is made for each individual connection, yielding a single pair of true-positive and false-alarm rates.

174 individual test can be made for each connection in the network, for example by comparing the estimated
175 value to a null distribution (bottom thread). Here again, we obtain two rates of false positives and
176 negatives.

177 *Coefficients from linear regression potentially predict better existing connections than residuals*

178 We start with the comparison between the predictability of coefficients and residuals for MV, GRu and
179 GRc for all connections in a given network. To do so, we simulate 500 randomly connected networks,
180 which are simulated with different sizes ($N = 50$ to 150 nodes), density and connectivity weights
181 (uniformly drawn in a randomly chosen range $[w_{\min}, w_{\max}]$); here inputs are *not* correlated: the ζ_i are
182 independent across node indices in Eq. (1). For each network configuration, we evaluate the accuracy for
183 connection detection via the area under the ROC curve (see the upper thread in Fig. 1B). Fig. 2A displays
184 this ROC-based accuracy as a function of the number of observed time samples (x-axis) represented by
185 violin plots for 500 randomly connected networks. When considering many samples (10^4), all methods
186 perform well. However, for smaller sample sets, the MVAR method exhibits superior performance than
187 GRu: as measured by the Mann-Whitney test, $p < 10^{-45}$, $p < 10^{-19}$ and $p < 10^{-5}$ for the three values of
188 observed samples, respectively.

189 Although error residuals ratios are in a different space from the true weights in A , one expects some
190 degree of correlation between them, such that Granger causality analysis effectively detects connections.
191 In Fig. 2B, both GRu and GRc estimates have a ranking similar to the original A weights (as measured by
192 the Spearman correlation) for $T = 10000$ observed time samples, but this weakens dramatically for
193 $T \leq 3000$. In contrast, the ranking for estimated MVAR coefficients reflects much better the original A
194 for $T \leq 3000$. In the studied networks, GRu performs slightly better than GRc. As analyzed in previous
195 studies, this can be consequence of the balance between redundant and synergistic activity exhibited by
196 the simulated network nodes (Stramaglia et al., 2014). To shed light into the effect of the network
197 structure, we next examine how the ROC-based performance in Fig. 2A depends on the controlled
198 network parameters. The four panels in Fig. 2C display the trends of the values for the 500 networks as a
199 function of the network size N , the network density, the minimum weight in the original network (w_{\min}
200 mentioned above) and the mean sum of incoming weights per node. For illustration purpose, the 500
201 networks are grouped in quartiles for each parameter. Not surprisingly, the estimation accuracy of all

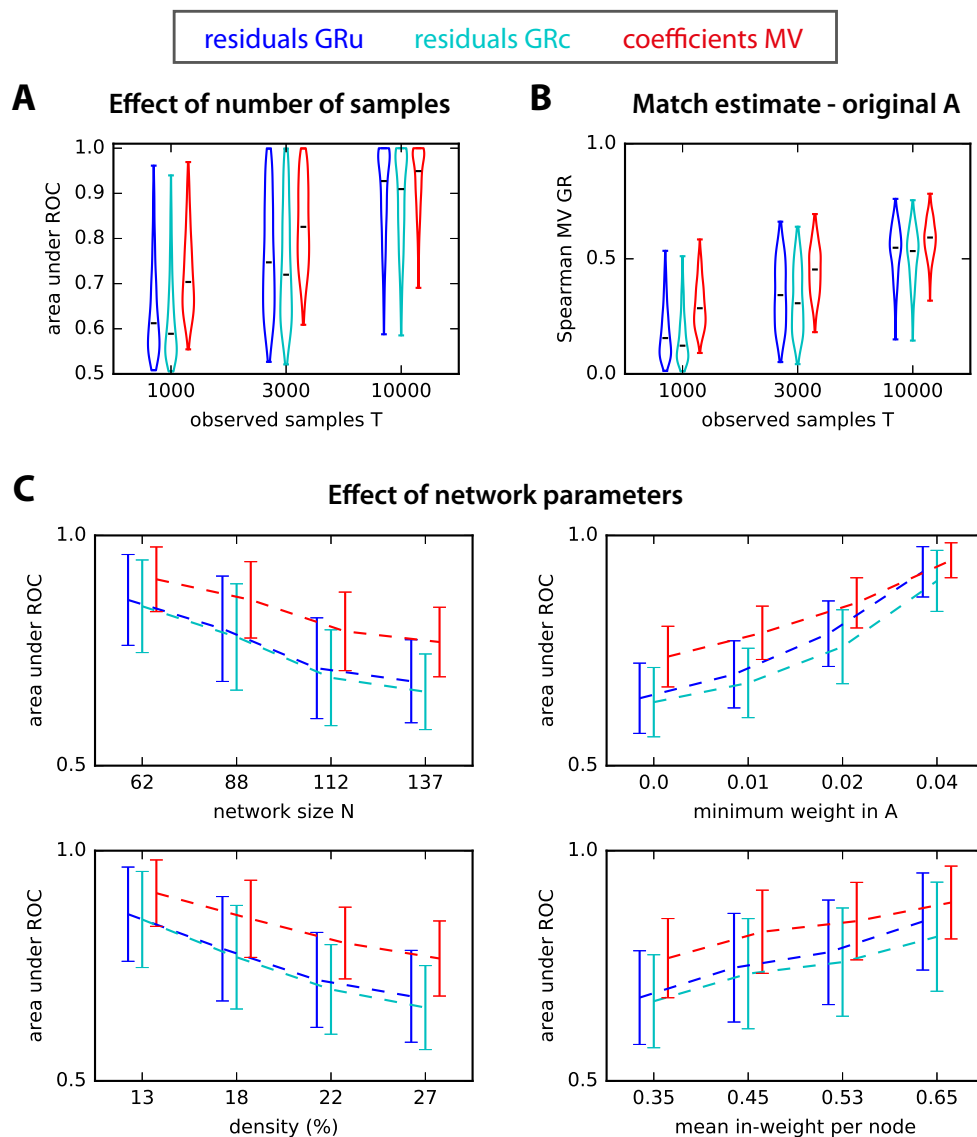
202 methods decreases as a function of the network size N and density, and increases as a function of the
203 minimum connectivity weight and the mean incoming weight per node. More interestingly, in
204 challenging configurations with small weights, MVAR consistently shows a superior performance by a
205 larger gap compared to Granger causality analyses. These findings support the use of coefficients to
206 robustly detect connections in recurrently connected networks. Note that GRu performs on average
207 slightly better than GRc here: the discrepancy decreases as a function of the network density, which may
208 follow from lower redundancy in the recurrent network (Stramaglia et al., 2014).

217 *A robust non-parametric significance test for MVAR*

218 We have so far examined the performance of different estimation methods based on the area under the
219 ROC curve, which corresponds to a single threshold for all connections in a network and combines the
220 information about false alarms and true detection over the whole range of estimated values. However, in
221 the context of real data, the decision for the existence of a connection typically relies on comparing the
222 value of the connection estimate with a given statistical threshold. For GRu and GRc, such parametric
223 tests have been developed, for example, based on the F statistics (Barnett & Seth, 2014). Equivalently, it
224 is sufficient to know how the values of the estimates for absent connections are distributed, in order to
225 select a desired rate of false alarms (type-1 error). In this section, we develop a significance test for the
226 estimated MVAR coefficients by providing a null-hypothesis distribution for absent connections.

227 Our approach relies on the fact that covariances reflect the underlying connectivity: we thus construct
228 the null distribution for estimates by performing a random permutation for each of the observed time
229 samples, which “destroys” the covariance structure apart from the variances on the diagonal of \widehat{Q}^0 as
230 illustrated in Fig. 3A; other methods will be tested in a later section. From the resulting covariance
231 matrices, we evaluate a surrogate connectivity matrix. The core result underlying our surrogate approach
232 is shown in Fig. 3B: the distribution of surrogate estimates (thick black line) is compared against the
233 distribution of existing (red) and absent (blue) connections in a simulated random network model: the
234 surrogate distribution in black provides a good approximation for the distribution of estimates for
235 non-existing connections in blue.

236 We consider two options - corresponding to the two threads in Fig. 1B - to test the existence of a
237 connection from an MVAR estimate while keeping the false alarm rate under control.

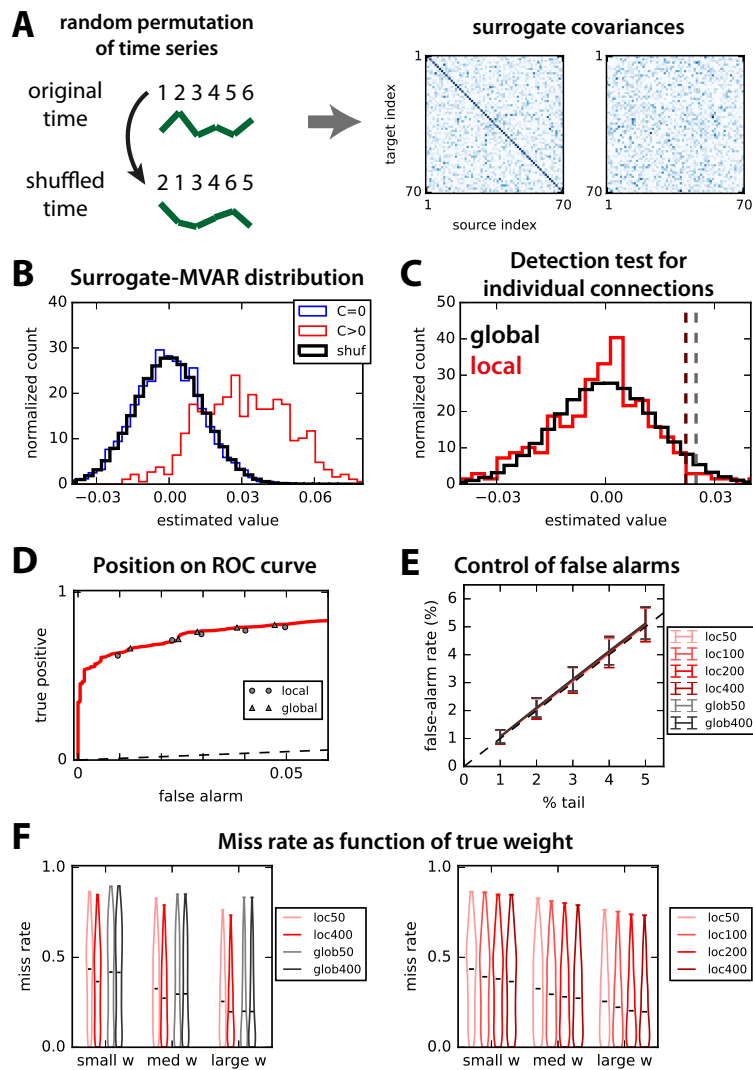


209 **Figure 2. ROC-based prediction power.** **A:** Area under ROC for estimated A obtain from log-ratios of residuals obtained from Granger causality analysis
 210 (unconditional for GRu and conditional for GRc) and MVAR. The x-axis indicates three sample size T for the observed network activity. The violin plots
 211 correspond to 500 simulated networks of various sizes and connectivity strengths (the horizontal black bar indicates the median). **B:** Match of the ranking
 212 between GRu, GRc and MVAR estimates and the original connectivity weights A , as measured by the Spearman correlation coefficient. The plotted values
 213 correspond to the 500 networks in **A** and the x-axis indicates the sample size T . **C:** Effect of network parameters on ROC-based performance. Influence of
 214 network size N , connectivity density, sum of recurrent connectivity strengths, minimum weight w_{\min} in A , mean noise on the diagonal of Σ and mean off-
 215 diagonal noise in Σ on the ROC-based accuracy in Fig. 2A. In each plot, the network configurations have been grouped in quartiles according to the parameter
 216 plotted on the x-axis, and the corresponding group mean and standard deviations are indicated; the curves are displaced horizontally to improve legibility.

- 238 ▪ The *global* test relies on the null distribution corresponding to the black histogram in Fig. 3C, which
239 is obtained by grouping together all SN^2 matrix elements of all matrices for $S = 200$ surrogates.
240 From that surrogate distribution, we perform a detection test by setting a threshold corresponding to
241 a percentage of the right tail equal to the desired false-alarm rate (here 2%), as illustrated by the
242 vertical gray dashed line.
- 243 ▪ Instead, the *local* test uses for each connection the surrogate distribution of S values, corresponding
244 to the same matrix element in each of the S surrogates. From that distribution in red in Fig. 3C, the
245 detection threshold is defined similarly (vertical dark red dashed line).

246 The rationale behind these two choices lies in the trade-off between taking into account spatial
247 heterogeneity in the network and gaining larger sample size, as illustrated by the distinct thresholds in
248 Fig. 3C. Note also that the F statistical test for Granger causality analysis corresponds to a global
249 threshold on the log ratio values. When varying the desired false-alarm rate, the two tests perform well,
250 as illustrated in Fig. 3D by their location close to the ROC curve (circles and triangles for local and
251 global, respectively).

252 To assess the effect of the small variability observed in Fig. 3D over the randomness of network
253 configurations, we simulate 500 randomly connected networks with the same parameters as in Fig. 2,
254 except for the size $50 \leq N \leq 90$ and the presence of input cross-correlations. Note that, from Fig. 2, the
255 chosen size N corresponds to a situation where Granger causality analysis performs relatively well as
256 compared to MVAR. The control of the false-alarm rate is displayed in Fig. 3E for both local and global
257 tests with various numbers S of surrogates. The control of false-alarm rates is close to perfect across
258 various values for all S and both tests (local and global), demonstrating the robustness of the proposed
259 method for randomly connected networks. Following, we fix the desired false-alarm rate to 2% and
260 evaluate the miss rate (true negatives) of both methods depending on the actual weight strength: in
261 Fig. 3F, connections are grouped in terciles for each network configuration. Interestingly, the local test
262 improves with the number S of surrogates (right panel), whereas the global test exhibits a constant
263 performance for all S (left panel). Note that the advantage of the local test over the global test
264 particularly concerns connections with small weights, which are difficult to detect, in line with Fig. 2C
265 (see influence of the minimum original weight).



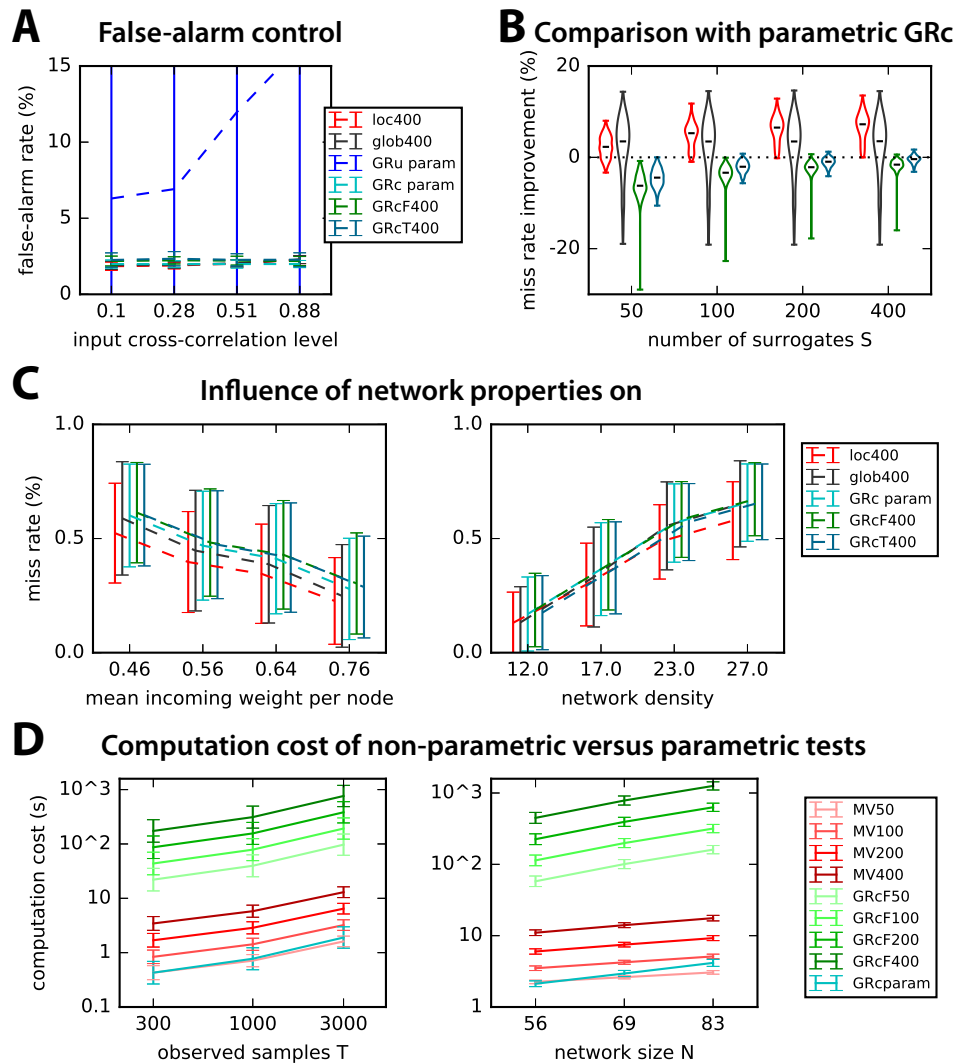
266 **Figure 3. Non-parametric test to assess significance for MVAR coefficients based on random permutations.** **A:** Schematic illustration of random
 267 permutation (numbers indicate time) applied independently to all observed time series (green curves) to generate the surrogate covariances (right panels). **B:**
 268 Pooled distributions of estimated weights for the existing (in red) and absent (blue) connections. The thick black curve indicates the distribution of connections
 269 over 100 surrogates, which closely matches the blue distribution. **C:** For a given connection, we compare two methods: for ‘local’ in red, the null distribution
 270 corresponds to the matrix elements for the same connection in 200 surrogates; for ‘global’ in black, the null distribution is the pooled distribution for all N^2
 271 elements of the 200 surrogate matrices (same as in B). The dark red and gray dashed lines indicate the detection thresholds corresponding to the 4% tail for
 272 those two options. **D:** The performance of the two non-parametric methods for the thresholds described in B and C is displayed on the ROC curve for a
 273 desired false-alarm rate ranging from 1 to 5%. Triangles indicate the local test and circles the global test. **E:** Comparison of the desired (% of the tail of null
 274 distribution) and actual rate of false alarms for the local and global tests when varying the the number S of surrogates (see figure in legend). Error bars indicate
 275 one standard deviation for 500 random networks; importantly, inputs for these networks have cross-correlation, unlike Fig. 2. **F:** Influence of the strength of
 276 original weight on the detection performance for the 500 random networks and a desired false-alarm rate set to 2% in E. In both panels, lighter colors indicate
 277 smaller numbers of surrogates S , in red for the local test and gray for the global test (see legends).

278 Because GRu does not take all network nodes into account, the presence of spatially correlated noise
279 (indicated by the purple dashed arrows in Fig. 1A) dramatically affects the false-alarm rate when using
280 the parametric significance F-test (Barnett & Seth, 2014), as shown in Fig. 4A by the dark blue dashed
281 curve. This is solved by the “complete” linear regression in GRc, achieving a quasi perfect control
282 irrespective of the input correlation level for both the parametric and two non-parametric tests (cyan,
283 green and blue-green dashed curves, respectively), as our non-parametric tests do (red and gray; recall
284 also Fig. 3E). We consider two non-parametric tests for GRc: ‘T’ stands for target, where the null
285 distribution of a connection is obtained by shuffling only the target, and ‘F’ stands for full, where we
286 shuffle all time series simultaneously as in our coefficient-based tests. Both perform equally in terms of
287 false alarms.

288 The main result of the paper is described in Fig. 4B, where the dashed line corresponds to the miss rate
289 for parametric GRc: our non-parametric method exhibits a better than miss rate - i.e., decrease - for both
290 local and global tests (in red and gray, respectively) on average over the same 500 random networks as in
291 Fig. 4A. For $S \geq 200$, the local test even becomes better in all cases. Note that the small miss-rate
292 improvements of about 7% actually corresponds to more than 50 existing connections per network here.
293 In contrast, both non-parametric tests for GRc perform worse than the parametric test here, with the
294 target-shuffling surrogate converging faster close to the non-parametric GRc. Fig. 4C displays the trends
295 of the performance of all 5 tests in Fig. 4B as a function of two network properties: the mean incoming
296 weight per node (left) and the density (right). The main result here is that the local test performs better
297 especially in difficult configurations with small weights and dense connectivity.

298 From Figs. 3F and 4B-C, we conclude that the local test is preferable to the global test provided
299 $S \geq 200$ surrogates are generated. However, the computational cost increases linearly with S , as
300 illustrated in Fig. 4D by the red curves. Note that the parametric GRc (in cyan) takes the same time to
301 calculate as $S = 50$ surrogates. However, our non-parametric method scales better than parametric GRc
302 when the network size increases. As a comparison, the full-network non-parametric test for GRc takes
303 longer time to compute, but further optimization of the calculations could be made that were not
304 incorporated here.

315 *Comparison of generation methods for surrogates for non-parametric MVAR*



305 **Figure 4. Comparison of our coefficient-based method with Granger causality analysis.** **A:** Comparison of the parametric tests for GRu (blue curve)
 306 and GRc (cyan) with the non-parametric methods for GRc (green for GRcF and blue-green for GRcT, see the text for details) and MVAR (red for local test and
 307 gray for global). The x-axis indicates the strength of input correlations (i.e., pink noise) in the simulated network. The desired false-alarm rate is set to 2% as
 308 in Fig. 3F and the number of observed time samples is $T = 3000$. Error bars indicate one standard deviation over the 500 random networks as in Fig. 3E. **B:**
 309 Comparison of the miss rate improvement (decrease) with respect to parametric GRc for the 500 networks in A as a function of the number S of surrogates
 310 (x-axis). Red indicates the local test, gray the global test, green the full-network non-parametric GRc and blue-green the target-only non-parametric GRc. **C:**
 311 Details of the performance of the 5 methods in B as a function of the mean incoming weight per node (left) and the network density (right). The plots for the
 312 miss rate are similar to those for the ROC-based prediction power in Fig. 2C. **D:** Comparison of the computational cost for the surrogate-based method and
 313 parametric tests as a function of the number T of observed samples (left) and network size (right). Only GRcF is shown, as GRcT takes much longer time in
 314 the unoptimized version that we use.

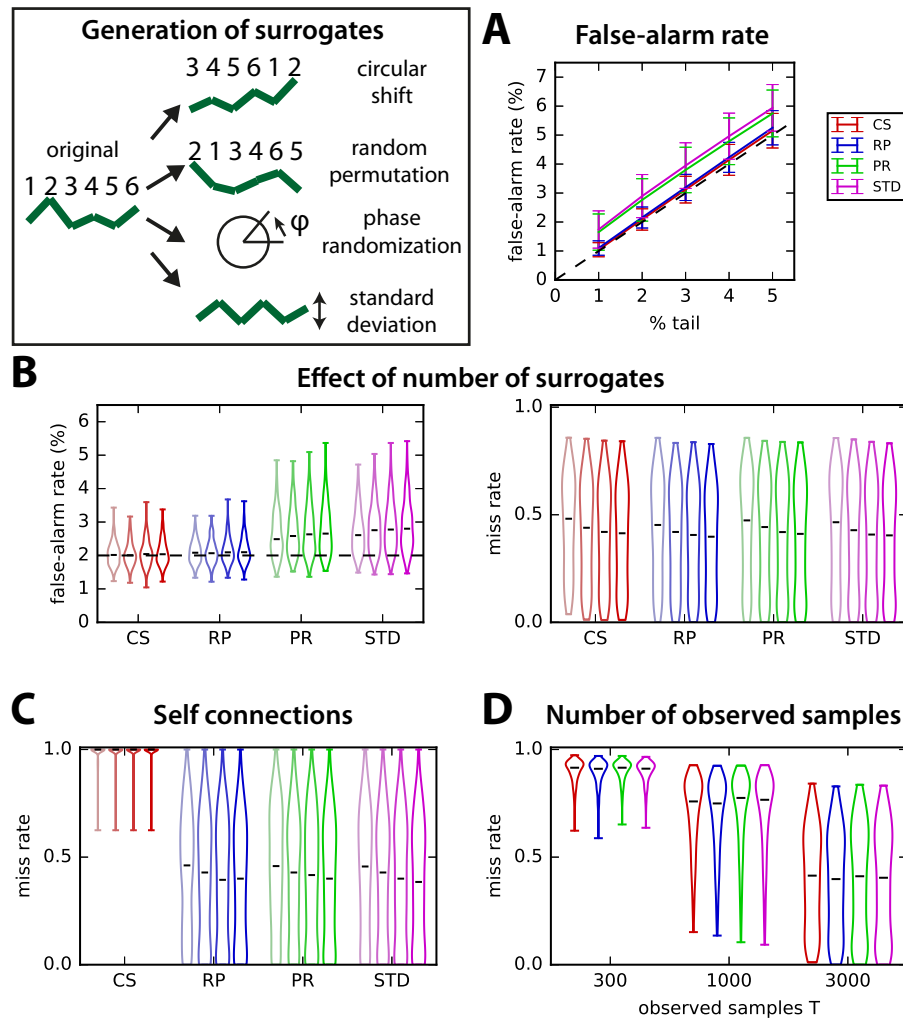
316 The fact that the OLS MVAR estimates can be obtained via the two covariance matrices (with and
317 without time shift, see Fig. 1A) hints at possible methods to generate surrogate by destroying the
318 information in these covariances. Methods to generate surrogate time series have been widely used in the
319 past: circular shifts of the time series (Faes et al., 2014), random permutation (Winkler et al., 2014) and
320 phase randomization (Schreiber & Schmitz, 1996) to generate a null distribution for the ratios in Eq. (4);
321 they are referred to here as CS, RP and PR, respectively. We thus consider these three methods (cf. box in
322 Fig. 5), as well as surrogate time series that only preserve the mean standard deviation averaged over the
323 network (STD), so as test to which extent it is important to preserve the spatial heterogeneity of the
324 nodes' activity. See Methods for details about the calculations.

325 The control of false alarms for local tests in Fig. 5A and B is better for CS and RP, whereas the
326 detection of true connections is similar for the four methods over 500 random networks of size $N = 70$.
327 However, CS fails to detect self connections (Fig. 5C). The reason is that, because CS surrogates preserve
328 the autocovariances in the time-shifted covariance, they fail to build a proper null distribution for
329 self-connections. The influence of the number of samples used in the estimation is similar for all
330 methods, as illustrated in Fig. 5D. The comparison with STD (purple), which averages the covariance
331 statistics over the whole network, suggests that the local test makes a good use of the heterogeneous
332 information across nodes. As a conclusion, we retain RP as the best option.

341 *Influence of network topology*

342 In this part, we test and compare the robustness of global and local surrogate-based detection tests to
343 specific connections and topological configurations. Here, $T = 3000$ observed samples and we compare
344 the local and global tests with $S = 400$ surrogates for 500 networks of each type. In all cases, the
345 simulated networks have the same size $N = 70$, but vary in connectivity density, distribution of recurrent
346 weights and level of input cross-correlation. We compare the miss rate for unidirectional, reciprocal and
347 self connections in the random networks examined until now (and a desired 2% of false alarms). Fig. 6A
348 shows that the miss rate is similar in unidirectional and reciprocal connections with the local test, which
349 performs slightly better than the global test (as in Fig. 3F).

350 Now we consider more elaborate network topologies than the random connectivity (Erdős-Rényi)
351 considered so far, namely modular and hierarchical networks. In Fig. 6B, we simulate 500 modular



333 **Figure 5. Comparison of the local test for 4 surrogate generation methods.** Surrogates are generated by independently performing for each original time
 334 series 1) random permutations, 2) circular shifts, 3) phase randomization and 4) replacing the original by a new random time series with the mean standard
 335 deviation averaged over all of the original time series in the network. **A:** Comparison of the control of the false-alarm rate for various thresholds on the tail of
 336 the distributions of 400 surrogates (% indicated on the x-axis). The error bars correspond to the variability over 500 random networks similar to Fig. 3 with
 337 $T = 3000$ observed time samples and the local test. **B:** Influence of the number S of surrogates on the detection performance for the 4 surrogate methods
 338 (x-axis). The violin plots indicate the distribution of false-alarm rates (left panel) and miss rates (right) for the 500 networks in A with a desired false-alarm
 339 rate set to 2% (dashed line in the left panel). Lighter to darker colors correspond to 50, 100, 200 and 400 surrogates, respectively. **C:** Same as the miss rate in
 340 B, but only for self connections. **D:** Influence of the number T of observed time samples on the miss rate for $S = 400$ surrogates.

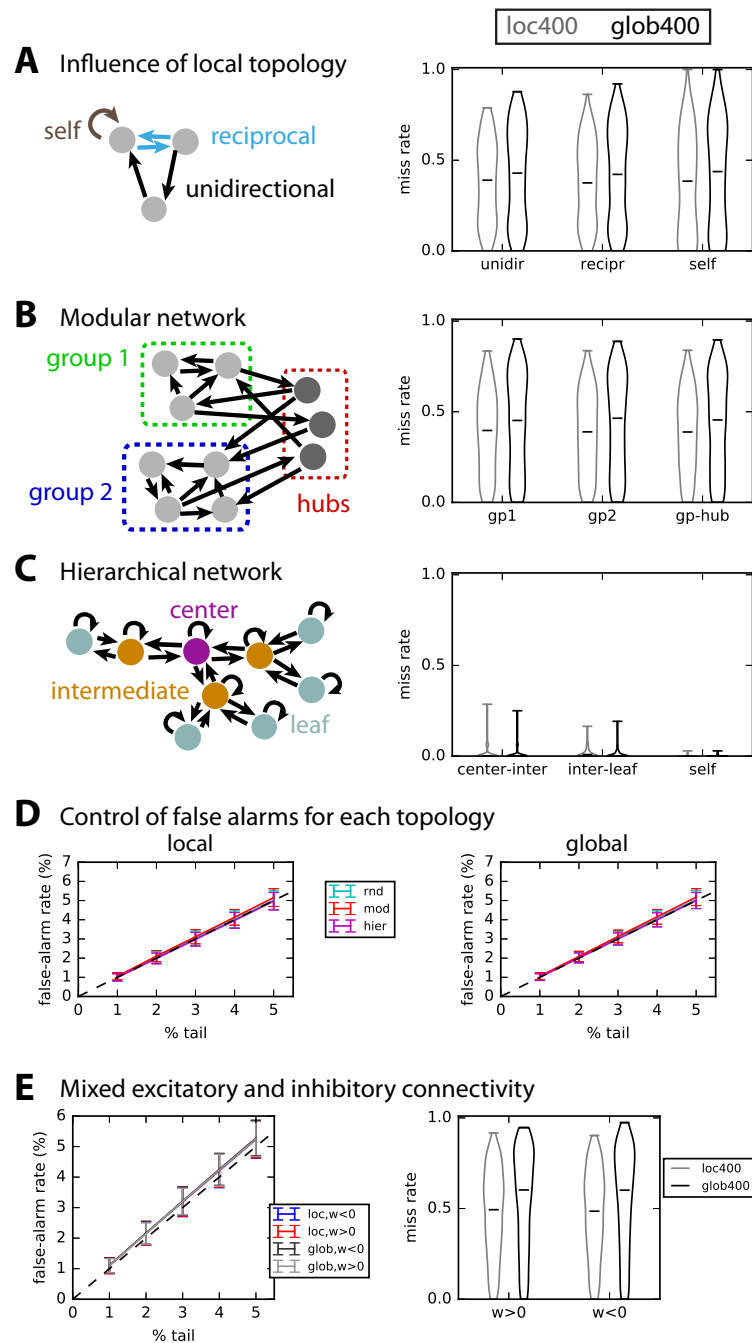
352 networks with two groups (green and blue) linked by hubs (red, about 5 to 15% of the nodes).
353 Interestingly, intra-group and hub-group connections have a similar miss rate with regard to using local
354 and global surrogates. In Fig. 6C, we simulate hierarchical networks of three layers, for which
355 connections either link the center and an intermediate node, or link an intermediate node and a leaf, or are
356 self connections. This network type is much sparser than the two types in A and B, yielding a quasi
357 perfect detection performance for all types of connections (miss rate < 0.1 in Fig. 6C). In all cases, the
358 local test performs better than the global test. However, the control of false-alarm rate is similar for both
359 tests with all topologies, as can be seen in Fig. 6D.

360 Finally, we consider a network with both excitatory and inhibitory connections (with a inhibitory ratio
361 equal to 5 to 50% of all) and perform the test by defining a threshold on both tails of the null distributions.
362 As can be seen in Fig. 6E, the positive/negative nature of the connection weights affect neither the
363 false-alarm nor the miss rate. However, the performance is poorer than with excitatory connections only.

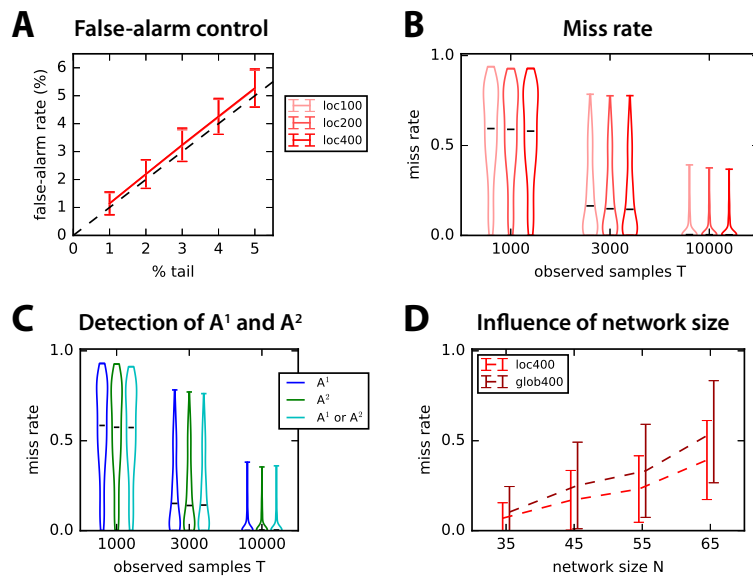
364 We conclude that, in those networks with spatial heterogeneity as with randomly connected networks,
365 the local test with an individual null distribution per connection performs better than the global test.
366 Recall that an improvement of the miss rate by 1% in a network with a density of 20% actually
367 corresponds to $N^2 0.2/100 \simeq 10$ existing connections here, so the plotted improvements concern about
368 50 connections.

376 *Applicability to second-order MVAR process*

377 As explained in Methods, an MVAR process whose state depends on the two previous time steps can be
378 estimated with the covariances with time shifts $\tau = 0, 1$ and 2; see Eqs. (11) and (12) for details. Here we
379 simply focus on random connectivity for the two corresponding matrices A^1 and A^2 , with size N that is
380 randomly drawn between 30 and 80; we construct A^1 and A^2 such that a connection $j \rightarrow i$ cannot be in
381 both matrices, but at most in one. The existing connections are detected with the non-parametric local
382 test relying on RP surrogates for each matrix separately, as a proof of concept. The control of false
383 alarms in Fig. 7A and the overall detection performance in Fig. 7B suggest that our surrogate method can
384 be extended satisfactorily to higher-order MVAR processes. Note that the improvement by generating
385 more surrogates is rather weak here. Importantly, there is no difference between the detection in A^1 and



369 **Figure 6. Robustness to non-trivial network topology.** **A:** Detection performance for unidirectional, reciprocal and self connections in 500 the randomly
 370 connected networks used so far in Fig. 3. **B:** Detection performance for modular topology schematically represented in the left: each of the 500 networks
 371 comprises of two groups connected by hubs. The connections are separated depending on whether the connect group nodes or hubs, as indicated by the
 372 diagram on the left. **C:** Similar to B with a hierarchical topology, where connections are grouped in 3 subsets: from center to intermediate nodes; from
 373 intermediate nodes to leaves; self connections. The results concern 500 such networks, which all have very low density. **D:** Control of false-alarm rate for the
 374 local (left) and global (right) significance tests and the three network topologies; the plot is similar to Fig. 3E. **E:** Control of false-alarm rate and miss rate for
 375 networks with both excitatory and inhibitory connections.



388 **Figure 7. Connectivity detection for second-order MVAR process.** **A:** Control of false-alarm rate for the local test with $S = 100, 200$ and 400 in
 389 both connectivity matrices A^1 and A^2 , corresponding to each time step. Error bars correspond to the variability over 500 network configurations with random
 390 connectivity and $T = 3000$ observed samples. **B:** Influence of the number T of observed samples (x-axis) on the miss rate for the 500 networks in A. Lighter
 391 to darker red indicates the number of surrogates S . **C:** Details of the detection performance for A^1 and A^2 separately, as well as connections in either A^1 or
 392 A^2 . The number of observed samples is indicated on the x-axis as in B. **D:** Influence of the network size N on the detection performance in C for the local and
 393 global tests with $S = 400$ surrogates.

386 A^2 , as demonstrated in Fig. 7C. Last, the network size worsens the miss rate in Fig. 7D, which affects
 387 more dramatically the global test as compared to the local test.

APPLICATION TO EXPERIMENTAL DATA

394 *Multiunit activity data obtained from Utah electrode array in monkey*

395 Now we consider data recorded from a monkey performing a visual task, where the stimulus corresponds
 396 to vertical gratings covering all recorded V1 receptive fields from the Utah arrays (see Methods for
 397 details). We aim to provide a proof of concept for the connectivity analysis for this type of data, so as to
 398 complement the more classical analysis based on the activity of individual channels; therefore we do not
 399 focus on comparing the 4 stimulus conditions with each other.

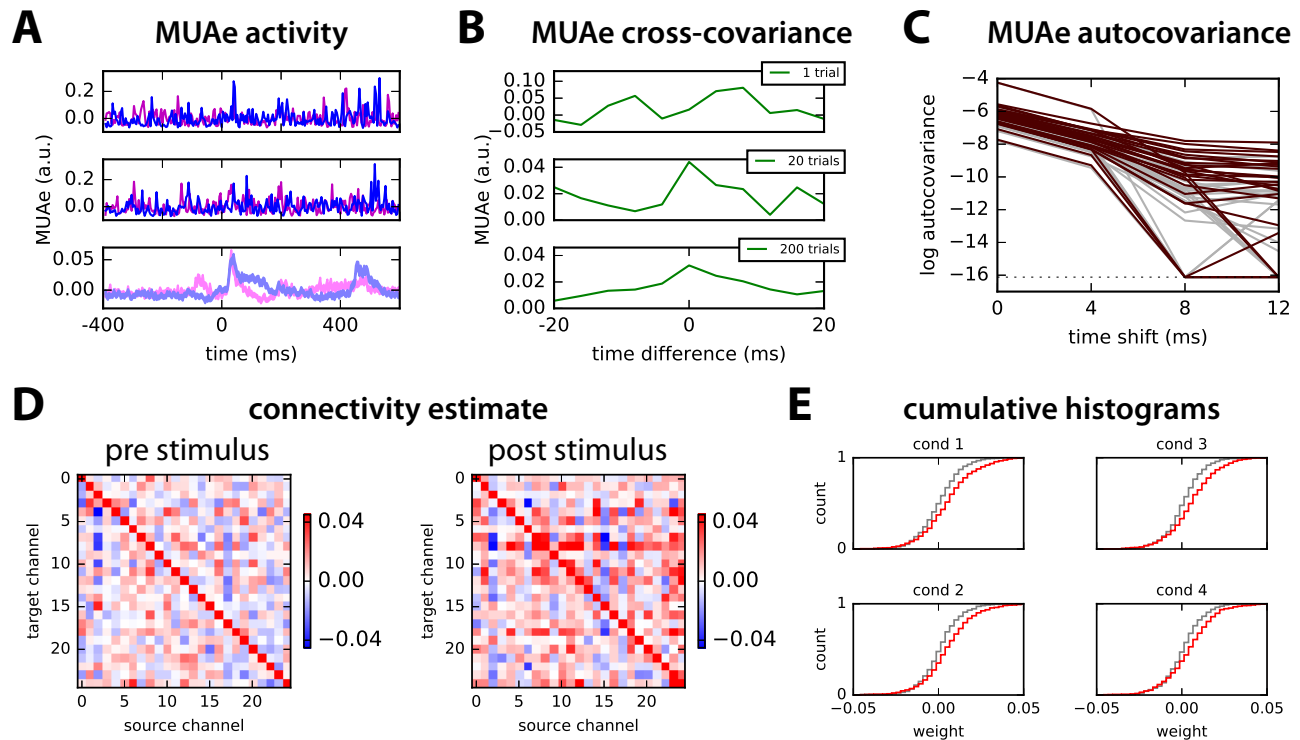
400 The multiunit activity envelope (MUAe) is obtained as described in Methods. In Fig. 8A, the resulting
401 MUAe is represented for two out of the 26 channels (red and purple) for two trials in the top and middle
402 panels, 400 ms before and 600 ms after the stimulus onset. The typical analysis of MUAe activity
403 consists of averaging over 200 trials, which exhibits a peak immediately after the stimulus for the two
404 channels in the bottom panel. Among the 26 channels, about a third show a large increase in activity after
405 the stimulus onset as compared to before (namely, a post-stimulus mean activity larger by more than three
406 standard deviations compared to the pre-stimulus activity); almost all channels show a moderate increase
407 of one standard deviation. One channel is discarded for a much larger activity (by 5 times) than all others.

408 To further investigate the temporal information conveyed by MUAe jointly for pairs of channels, we
409 calculate the pairwise covariances between them, after centering the MUAe activity individually for each
410 trial. Fig. 8B shows the stabilization of the cross-covariance between the two channels in Fig. 8A from a
411 single trial to averages over 20 and 200 trials. Note the asymmetry with respect to time difference: this
412 information is extracted by the network model to estimate the interactions between the neuronal
413 populations recorded by the channels. Then we verify that the model can be applied to these data, by
414 examining the MUAe autocovariances in Fig. 8C, which exhibit a profile corresponding to an exponential
415 decay up to two time shifts (i.e., 8 ms for the downsampling every 4 ms), that is, a straight line in the log
416 plot. This suits an autoregressive model with large positive values on the diagonal of the connectivity
417 matrix A .

418 Both connectivity matrices for the 25 channels estimated using the MVAR method before and after the
419 stimulus are illustrated in Fig. 8D for condition 1: we find larger off-diagonal values for the period after
420 the stimulus than before. This is actually true for all conditions, as indicated by the more spread
421 distributions in red as compared to gray in Fig. 8E. The channels appear to be coordinated at the
422 considered time scale of 4 ms and their collective interaction scheme is affected by the stimulus
423 presentation.

433 *Significance test for real data: interactions related to stimulus presentation*

434 We then use the local and global tests based on 1000 surrogates (with random permutation) to retain only
435 significant interactions from the estimates in Fig. 8D: this leaves a few interactions for the pre-stimulus
436 period in Fig. 9A (left panel), 8 out of 650, which is of the order of the desired false-alarm rate set to 1%



424 **Figure 8. Application to multiunit activity (MUAe) data.** **A:** Example of two trials (top and middle panels) of multiunit activity envelope (MUAe) for
 425 two channels of recordings using Utah electrode array in the primary visual cortex of a monkey (in arbitrary units; see text for further details). The bottom
 426 panel represents the average over 200 trials (with standard-error mean for the thickness of the curve). The stimulus is presented at time point 100 (actually 400
 427 ms, since the smoothing corresponds to a smoothing window of 4 ms). **B:** Example of cross-covariance between the two channels in A averaged over 1, 20
 428 and 200 trials. **C:** Autocovariance profiles of MUAe signals for all 25 channels and time shifts up to 12 ms averaged over 200 trials plotted with a log y-axis:
 429 comparison of signals before (gray) and after (red) the stimulus presentation. **D:** MVAR estimates of the connectivity between the 25 channels for the MUAe
 430 activity 200 ms before and after stimulus presentation (i.e., 50 time points each), averaged over 200 trials. The scaling has been optimized to enhance the
 431 legibility off-diagonal elements. **E:** Comparison of the cumulative distribution of connectivity weights (off-diagonal elements in D) for the 4 conditions. Gray
 432 and red indicate before and after the stimulus, respectively.

437 (namely, the extreme 0.5% of each tail). In contrast, many more post-stimulus interactions survive the
438 significance tests in the right panel: almost all these interactions are unidirectional. The counterpart for
439 circular shift for Fig. 9B involve 24 interactions in common with Fig. 9A. On average over the 4
440 conditions, 22 post-stimulus interactions are common between the two shuffling methods, to be
441 compared with 7 for the pre-stimulus period (both with a standard deviation of 4); this corresponds to
442 3.5% of all possible interactions. Almost all detected interactions are unidirectional, as illustrated in
443 Fig. 9C for both local and global tests for the post-stimulus period. Varying the threshold on the tail of
444 the null distributions, we see that the number of detected interactions is close to the desired false-alarm
445 rate for the pre-stimulus period in Fig. 9D (dark red and black curves, respectively). In contrast,
446 post-stimulus interactions are many more for both local and global tests (light red and gray). The global
447 test detects fewer interactions than the local test, indicating the necessity to take into account the
448 disparities across channels. Around 57% of post-stimulus interactions detected by the global test (largest
449 values in absolute value) are found by the local test.

450 Finally, we check the relationship between the strengths of significant interactions - in absolute value -
451 and the increase of average MUAe observed in Fig. 8A (lower panel). In Fig. 9E, the plotted dots
452 correspond to the pre-post change in the sum of incoming (left panel) and outgoing (right panel)
453 significant interactions for each node. The summed interaction values positively correlate with the MUAe
454 difference (post minus pre) only for the incoming interactions: $p = 0.03$ with a coefficient of 0.21;
455 nevertheless, the plotted values exhibit a large variability, which moderates this significance. In contrast,
456 outgoing interactions exhibit a non-significant negative correlation ($p \gg 0.1$). This suggests a
457 stimulus-driven gating of the effective gain for incoming anatomical connections to recorded cell
458 populations. The application of our method thus unravels stimulus-driven directed interactions and
459 cannot be merely explained by an increase of single-channel MUAe activity.

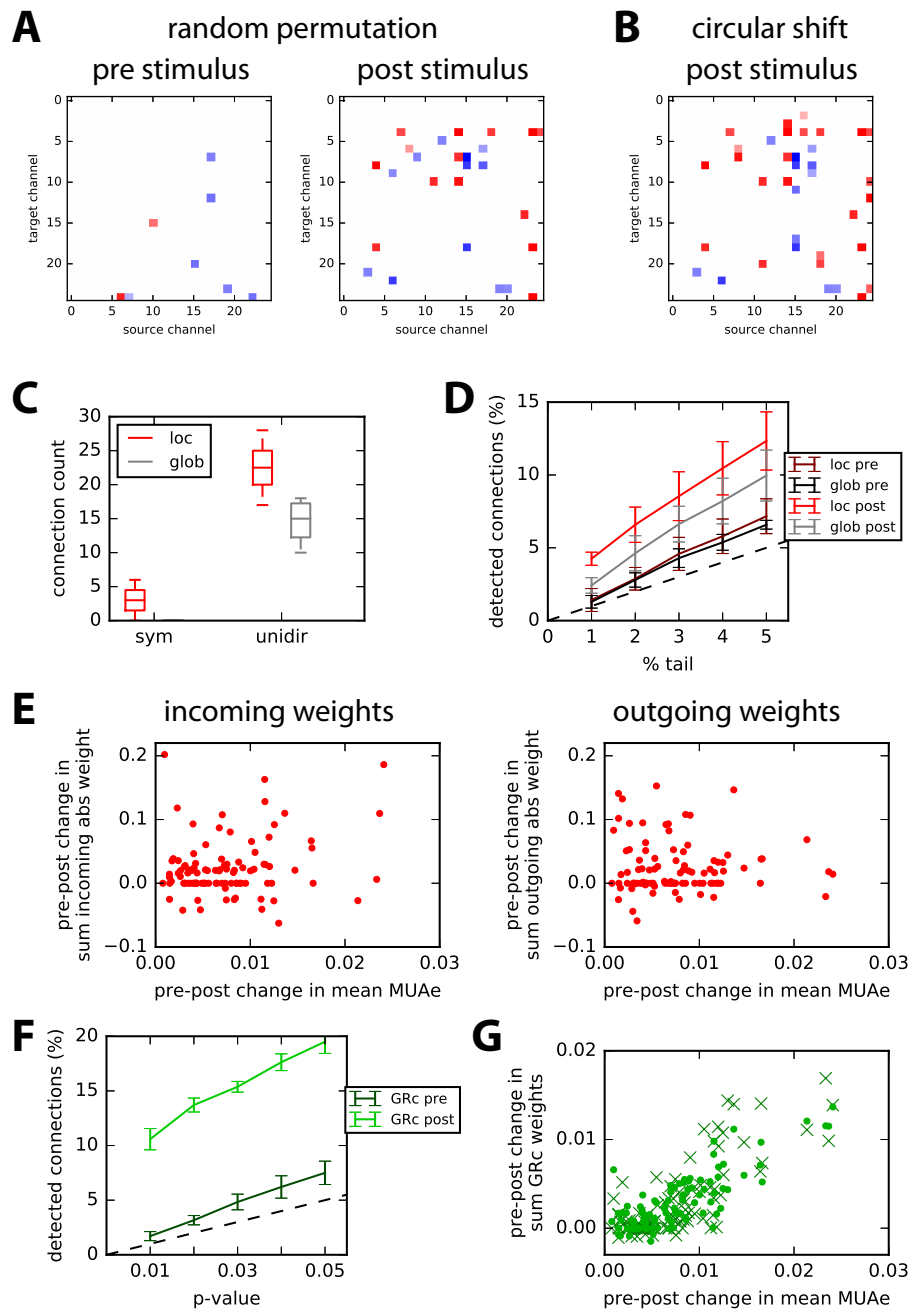
460 In comparison, similar detection with parametric GRc testing gives more interactions for the
461 post-stimulus period in Fig. 9F (bright green), more than twice the number for a p-value of 0.01
462 (corresponding to 1% in Fig. 9D). Moreover, the distributions of MVAR coefficients in Fig. 9E and the
463 corresponding ones for GRc estimated values have similar KS distance when comparing - for each
464 condition - the pre- and post-periods (mean of 0.17 with a std of 0.01 over the 4 conditions). This means
465 that GRc values collectively discriminate between the two periods as well as the estimated MVAR

466 coefficients. However, the pre-post changes in incoming and outgoing GRc values strongly correlate with
467 the change in mean MUAe ($p < 10^{-20}$ for both incoming and outgoing connections); note that the
468 estimated connectivity is not symmetric, though. This contrasts with the two plots in Fig. 9E and
469 suggests that the two methods may capture distinct effects at work in the network of neuronal
470 populations. Note that the non-parametric method for GRc did not work for the experimental data here,
471 failing to detect more interactions than the expected false-alarm rate.

DISCUSSION

481 *Non-parametric MVAR-based detection of linear feedback in recurrent networks*

482 This paper proposed a non-parametric method to detect pairwise feedback connections in biological
483 networks with possibly strong and/or dense recurrent feedback. We examined the benefit of detecting
484 directional connections in MVAR-like models estimated using the OLS autoregressive coefficients
485 instead of the error residuals ratios (Fig. 1). To our understanding, the good performance of the presented
486 method has three reasons. First, the ROC-based prediction power in Fig. 2, which relies on the estimated
487 ranking of connections in the network (i.e., from small to strong weights), is more robust for the
488 regression coefficients than residual log ratios for recurrent networks with relatively large density
489 (0.1 – 0.3%); these networks overall imply many redundancy and convergence patterns of connections
490 (M. Ding et al., 2006; Stramaglia et al., 2014). Second, practical coefficient-based connectivity detection
491 performed using non-parametric significance tests based on time-series randomization (red distribution in
492 Fig. 4B) yields better results than conditional Granger causality ratio test in the time domain, either with
493 the standard parametric F test (Barnett and Seth (2014), dashed line) or with non-parametric testing
494 (green and blue-green distributions). Finally, the use of connection-specific significance testing achieves
495 higher accuracy than a network-pooled alternative, especially when asymptotic assumptions do not hold
496 (e.g., small number of time samples), as illustrated by the local non-parametric test in Figs. 3F and 4B
497 (red versus dark gray). Together, our results highlight the need for testing strategies that capture the
498 heterogeneity of sufficiently large networks to detect individual connections. Further note that assessing
499 the connectivity via the regression coefficients space brings an additional advantage for network studies:
500 the estimated connection weights can be interpreted and compared across the whole network, for
501 example using graph theory (Sporns, 2013).



472 **Figure 9. Detection of significant interactions.** **A:** Examples of significant interactions with the time-rolling surrogate method; left and right panels
 473 correspond to pre- and post-stimulus periods, respectively. The p-value corresponds to the upper and lower 0.5% tails of 1000 surrogates (local test). Many
 474 more interactions are found for post- than pre-stimulus period. **B:** Matrix of post-stimulus interactions where green pixels indicate symmetric and purple
 475 asymmetric significant interactions; other are left blank. **C:** Comparison of number of asymmetric interactions versus symmetric interactions for the local (red)
 476 and global (gray) tests over the 4 conditions. **D:** Ratios of detected interactions for pre- and post-stimulus periods with the desired false-alarm rate equal to
 477 1 to 5% corresponding to both local and global tests. **E:** Change between pre- and post-stimulus periods of the sum of incoming (left) and outgoing (right)
 478 significant weights - in absolute value - plotted against the change in the change in mean MUAe. Each dot represents a channel and all 4 conditions are grouped
 479 here. **F:** Similar plot to D for parametric GRC. **G:** Similar plot to E with sums of GRC values for each node over its incoming (dots) and outgoing (crosses)
 480 connections. Only significant interactions passing the parametric test with p-value of 0.01 in F are retained here.

502 Our approach for generating surrogate distributions can be encompassed in the family of constrained
503 randomization methods (Schreiber, 1998; Schreiber & Schmitz, 1996). Here we have shown that random
504 permutation provides a good estimation of all types of connections (Fig. 5) in false-alarm and miss rates.
505 Comparatively, the circular-shift method performs as well except for self connections that are not
506 detected at all; this also holds for non-random topologies (results not shown). Hence, preserving the
507 autocovariance structure in the generation of surrogates does not provide a substantial advantage here
508 (Fig. 5B-C). Both methods show a good control for the false-alarm rate in comparison to phase
509 randomization (Schreiber & Schmitz, 1996) and a control Gaussian approximation over the whole
510 network (STD), which lead to an excess of about 1% of false alarms (i.e., ~ 50 connections for a network
511 of 70 nodes). These results show the importance of choosing a surrogate method adapted to the detection
512 problem. For distinct dynamics governing the nodal activity such as nonlinearities, conclusions may
513 differ and further research along these lines is necessary.

514 As mentioned earlier, the use of an individual null distribution for each connection (local test) gives
515 better results for the miss rate (by a few %) than lumping together all matrix elements of all surrogates
516 (global test), provided sufficiently many surrogates are generated. For the size of networks considered
517 here, computation time is not an issue (Fig. 4D) and our results support the choice of the local test over
518 the global test to attain between accuracy in the true-positive detection. This may be especially true for
519 specific topologies or networks with both excitatory and inhibitory connections, see Fig. 6. In other
520 words, the local test incorporates to a better extent the network heterogeneities in order to build the null
521 distribution for each connection. The present study was limited to ordinary least-square (OLS) estimates
522 for MVAR, but there exist alternative estimators such as the locally weighted least-square regression
523 (Ruppert & Wand, 1994) that may perform better for particular network topologies. The extension of the
524 presented surrogate techniques to the case where observations are sparser than connections - implying
525 that the covariance matrix is not invertible - is another interesting direction to explore (Castelo &
526 Roverato, 2006).

527 The problem of multiple comparison is intrinsic to brain connectivity detection as the number of
528 testable connections across brain regions is massive (Rubinov & Sporns, 2010). In this context, different
529 approaches have been developed to control the family-wise error rate in the weak sense. For instance,
530 many studies on neuroimaging data (Genovese, Lazar, & Nichols, 2002; Nichols & Hayasaka, 2003) or

531 electrophysiology (Lage-Castellanos, Martínez-Montes, Hernández-Cabrera, & Galán, 2010) have
532 resorted to procedures that control the false discovery rate (FDR) (Benjamini & Hochberg, 1995), namely,
533 the expected number of falsely declared connections among the total number of detections. These
534 methods make decisions on single connections relying on the entire sequence of p-values computed for
535 each connection and yield substantial statistical power gains over more conservative methods such as
536 Šidák-Bonferroni (Abdi, 2007). With the ever growing application of graph theory to brain connectivity,
537 new methods have been proposed that exploit the clustered structure of the the declared connections
538 (Han, Yoo, Seo, Na, & Seong, 2013; Zalesky, Fornito, & Bullmore, 2010) to propose cluster-based
539 statistical tests (Maris & Oostenveld, 2007) that attain similar performance to FDR methods. The present
540 work can therefore be understood as a primary step before performing any or several multiple-correction
541 procedures. By defining an accurate null model of inexistent connections, p-value estimates per
542 connection are improved and cluster-based surrogate distributions can be better approximated, which is
543 expected to empower the overall control of false positive rates in network connectivity analysis.

544 *Applications to real electrophysiological and neuroimaging data*

545 A motivation for our method is the detection of neuronal interactions between electrode channels, which
546 is often performed using spectral Granger causality analysis on local-field potential (low-passed signal of
547 electrode measurements) or ECOG measurements. As an alternative, we have applied our connectivity
548 detection method to MUAe recorded from macaque area V1 in order to provide a proof of concept.
549 Multi-channel recording devices have been developed in the past years to obtain this type of data (Fan et
550 al., 2011; Roy & Wang, 2012) and a recent cognitive study has highlighted group properties of MUAe
551 activity for a similar experimental setup to the one used here Engel et al. (2016). Looking at the
552 variability for individual trials in Fig. 8A (upper and middle panels), it is rather surprising that the MUAe
553 conveys temporal information about joint activity for pairs of channels (Fig. 8B), which can be related
554 to causal directed interactions. This means that the high trial-to-trial variability is not an absolute
555 limitation to temporal coordination, even though the latter only becomes apparent over multiple trial
556 repetitions, just as the post-stimulus increase of MUAe for single channels (lower panel in Fig. 8A). The
557 procedure detects a number of significant directional interactions well above the false positive rate that
558 were not merely explained by the MUAe increase (Fig. 9D and E). In comparison, the control for the
559 pre-stimulus period in the 4 different conditions detects just above the expected false-alarm rate. We find

560 that the parametric Granger causality test also detects many more interactions after the stimulus than
561 before; however, these interactions happen to link channels exhibiting the strongest increases in MUAe
562 activity. This presents the caveat of providing little information in addition to the changes observed at the
563 single-node level, when interpreting the obtained results. The research of optimal preprocessing - in
564 particular the filtering to obtain MUAe - to obtain a robust detection of interactions is left for a later
565 study. Likewise, such electrode recordings are often analyzed with respect to specific frequency bands
566 (e.g., alpha and gamma), but the adaptation of our framework along the lines of previous works (Dhamala
567 et al., 2008; L. Ding et al., 2007) and comparison of detected interactions with established methods will
568 be done in future research.

569 More generally, our methodology requires adequate preprocessing of multivariate time series - activity
570 aggregation over hundreds of voxels for fMRI and 4-ms smoothing of MUAe for electrode recordings -
571 such that the autocovariance profiles match the exponential decay of the dynamic network model with
572 linear feedback, which underlies the connectivity analysis. Although filtered and smoothed signals fall
573 into the class of autoregressive moving-average (ARMA) models, our approach is applicable provided
574 the autocovariances exhibit a profile resembling Fig. 8C. We further expect non-parametric testing
575 methods for be extendable to ARMA processes, complementing approaches developed by Barnett and
576 Seth (2015); Friston et al. (2014). In theory, stationarity of the time series remains a critical issue as we
577 need sufficiently many observed samples to obtain precise covariances from which we estimate the
578 connectivity, but this may not be a strong limitation for MUAe in practice in view of the post-stimulus
579 average response shown in Fig. 8A.

ACKNOWLEDGMENTS

580 MG acknowledges funding from the Marie Skłodowska-Curie Action (grant H2020-MSCA-656547).
581 MG and GD were supported by the Human Brain Project (grant FP7-FET-ICT-604102 and
582 H2020-720270 HBP SGA1). GD and ATC were supported by the European Research Council Advanced
583 Grant DYSTRUCTURE (Grant 295129). AT was supported by the UK Medical Research Council (grant
584 MRC G0700976). The authors are grateful to Robert Castelo and Inma Tur for constructive discussions.

AUTHOR CONTRIBUTIONS

585 Project was formulated by MG and ATC. Simulation code was developed by MG. Experimental data
586 were provided by AT and XC. Manuscript was written by MG, ATC, AT and GD.

587

588

REFERENCES

589

590 Abdi, H. (2007). The bonferonni and šidák corrections for multiple comparisons. *Encyclopedia of measurement and*
591 *statistics*, 3, 103-107.

592 Amemiya, T. (1974). Multivariate regression and simultaneous equation models when the dependent variables are truncated
593 normal. *Econometrica: Journal of the Econometric Society*, 999-1012.

594 Babiloni, F., Cincotti, F., Babiloni, C., Carducci, F., Mattia, D., Astolfi, L., . . . He, B. (2005). Estimation of the cortical
595 functional connectivity with the multimodal integration of high-resolution eeg and fmri data by directed transfer function.
596 *Neuroimage*, 24(1), 118-131.

597 Barnett, L., & Seth, A. (2014). The MVGC multivariate granger causality toolbox: A new approach to granger-causal
598 inference. *J Neurosci Methods*, 223, 50-68.

599 Barnett, L., & Seth, A. (2015). Granger causality for state-space models. *Phys Rev E*, 91, 040101. doi:
600 10.1103/PhysRevE.91.040101

601 Barrett, A., & Barnett, L. (2013). Granger causality is designed to measure effect, not mechanism. *Front Neuroinform*, 7, 6.
602 doi: 10.3389/fninf.2013.00006

603 Benjamini, Y., & Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple
604 testing. *Journal of the royal statistical society. Series B (Methodological)*, 289-300.

605 Castelo, R., & Roverato, A. (2006). A robust procedure for Gaussian graphical model search from microarray data with p
606 larger than n. *The Journal of Machine Learning Research*, 7, 2621-2650. Retrieved from
607 <http://dl.acm.org/citation.cfm?id=1248641>

608 Dhamala, M., Rangarajan, G., & Ding, M. (2008). Analyzing information flow in brain networks with nonparametric
609 granger causality. *Neuroimage*, 41(2), 354-362.

610 Diks, C., & DeGoede, J. (2001). A general nonparametric bootstrap test for granger causality. In (p. 391-403). Bristol, UK:
611 Institute of Physics Publishing.

- 612 Ding, L., Worrell, G., Lagerlund, T., & He, B. (2007). Ictal source analysis: localization and imaging of causal interactions
613 in humans. *Neuroimage*, *34*(2), 575-586.
- 614 Ding, M., Chen, Y., & Bressler, S. (2006). Granger causality: basic theory and application to neuroscience. In B. Schelter,
615 M. Winterhalder, & J. Timmer (Eds.), *Handbook of time series analysis: Recent theoretical developments and*
616 *applications*. arXiv preprint q-bio/0608035: Wiley-VCH Verlage.
- 617 Engel, T. A., Steinmetz, N. A., Gieselmann, M. A., Thiele, A., Moore, T., & Boahen, K. (2016). Selective modulation of
618 cortical state during spatial attention. *Science*, *354*, 1140-1144.
- 619 Faes, L., Marinazzo, D., Montalto, A., & Nollo, G. (2014). Lag-specific transfer entropy as a tool to assess cardiovascular
620 and cardiorespiratory information transfer. *IEEE Trans Biomed Eng*, *61*, 2556-68. doi: 10.1109/TBME.2014.2323131
- 621 Fan, D., Rich, D., Holtzman, T., Ruther, P., Dalley, J. W., Lopez, A., ... Yin, H. H. (2011). A wireless multi-channel
622 recording system for freely behaving mice and rats. *PLoS One*, *6*, e22033.
- 623 Friston, K., Bastos, A., Oswal, A., van Wijk, B., Richter, C., & Litvak, V. (2014). Granger causality revisited. *Neuroimage*,
624 *101*, 796-808. doi: 10.1016/j.neuroimage.2014.06.062
- 625 Friston, K., Harrison, L., & Penny, W. (2003). Dynamic causal modelling. *Neuroimage*, *19*(4), 1273-1302.
- 626 Genovese, C., Lazar, N., & Nichols, T. (2002). Thresholding of statistical maps in functional neuroimaging using the false
627 discovery rate. *Neuroimage*, *15*(4), 870-878.
- 628 Geweke, J. (1982). Measurement of linear dependence and feedback between multiple time series. *J Am Stat Assoc*, *77*,
629 304-313.
- 630 Geweke, J. (1984). Measures of conditional linear dependence and feedback between time series. *J Am Stat Assoc*, *79*,
631 907-915.
- 632 Granger, C. (1969). Investigating causal relations by econometric models and cross-spectral methods. *Econometrica:*
633 *Journal of the Econometric Society*, 424-438.
- 634 Han, C., Yoo, S., Seo, S., Na, D., & Seong, J.-K. (2013). Cluster-based statistics for brain connectivity in correlation with
635 behavioral measures. *PLoS one*, *8*(8), e72332.

- 636 Kaminski, M., & Blinowska, K. (1991). A new method of the description of the information flow in the brain structures.
637 *Biological cybernetics*, 65(3), 203-210.
- 638 Kamiński, M., Ding, M., Truccolo, W., & Bressler, S. (2001). Evaluating causal relations in neural systems: Granger
639 causality, directed transfer function and statistical assessment of significance. *Biological cybernetics*, 85(2), 145-157.
- 640 Lage-Castellanos, A., Martínez-Montes, E., Hernández-Cabrera, J., & Galán, L. (2010). False discovery rate and
641 permutation test: an evaluation in erp data analysis. *Statistics in medicine*, 29(1), 63-74.
- 642 Li, Y., Ye, X., Liu, Q., Mao, J., Liang, P., Xu, J., & Zhang, P. (2016). Localization of epileptogenic zone based on graph
643 analysis of stereo-eeg. *Epilepsy Res*, 128, 149-157.
- 644 Lusch, B., Maia, P., & Kutz, J. (2016). Physical review e covering statistical, nonlinear, biological, and soft matter physics
645 highlights recent accepted authors referees search press about inferring connectivity in networked dynamical systems:
646 Challenges using granger causality. *Phys Rev E*, 94, 032220.
- 647 Lütkepohl, H. (2005). *New introduction to multiple time series analysis*. Springer Science & Business Media.
- 648 Maris, E., & Oostenveld, R. (2007). Nonparametric statistical testing of eeg-and meg-data. *Journal of neuroscience*
649 *methods*, 164(1), 177-190.
- 650 Massey, J. (1990). Causality, feedback and directed information. In *Proc. int. symp. inf. theory applic.(isita-90)*
651 (p. 303-305).
- 652 Messé, A., Rudrauf, D., Benali, H., & Marrelec, G. (2014). Relating structure and function in the human brain: relative
653 contributions of anatomy, stationary dynamics, and non-stationarities. *PLoS Comput Biol*, 10, e1003530.
- 654 Michalareas, G., Schoffelen, J., Paterson, G., & Gross, J. (2013). Investigating causality between interacting brain areas
655 with multivariate autoregressive models of meg sensor data. *Human Brain Mapping*, 34(4), 890-913. Retrieved from
656 <http://dx.doi.org/10.1002/hbm.21482> doi: 10.1002/hbm.21482
- 657 Nedungadi, A., Rangarajan, G., Jain, N., & Ding, M. (2009). Analyzing multiple spike trains with nonparametric granger
658 causality. *J Comp Neu*, 27(1), 55-64.
- 659 Nichols, T., & Hayasaka, S. (2003). Controlling the familywise error rate in functional neuroimaging: a comparative review.
660 *Statistical methods in medical research*, 12(5), 419-446.

- 661 Rogers, B., Katwal, S., Morgan, V., Asplund, C., & Gore, J. (2010). Functional mri and multivariate autoregressive models.
662 *Magnetic Resonance Imaging*, 28(8), 1058-1065.
- 663 Roy, S., & Wang, X. (2012). Wireless multi-channel single unit recording in freely moving and vocalizing primates. *J*
664 *Neurosci Methods*, 203, 28-40.
- 665 Rubinov, M., & Sporns, O. (2010). Complex network measures of brain connectivity: uses and interpretations. *Neuroimage*,
666 52(3), 1059-1069.
- 667 Ruppert, D., & Wand, M. (1994). Multivariate locally weighted least squares regression. *Ann Statist*, 22, 1346-1370.
- 668 Schreiber, T. (1998). Constrained randomization of time series data. *Phys Rev Lett*, 80, 2105.
- 669 Schreiber, T. (2000). Measuring information transfer. *Phys Rev Lett*, 85, 461.
- 670 Schreiber, T., & Schmitz, A. (1996). Improved surrogate data for nonlinearity tests. *Phys Rev Lett*, 77, 635.
- 671 Seth, A., Barrett, A., & Barnett, L. (2015). Granger causality analysis in neuroscience and neuroimaging. *J Neurosci*, 35,
672 3293-3297. doi: 10.1523/JNEUROSCI.4399-14.2015
- 673 So, K., Koralek, A., Ganguly, K., Gastpar, M., & Carmena, J. (2012). Assessing functional connectivity of neural ensembles
674 using directed information. *J Neu Eng*, 9(2), 026004.
- 675 Sporns, O. (2013). Making sense of brain network data. *Nat Methods*, 10, 491-493.
- 676 Storkey, A., Simonotto, E., Whalley, H., Lawrie, S., Murray, L., & Mcgonigle, D. (2007). Learning structural equation
677 models for fmri. In B. Schölkopf, J. Platt, & T. Hoffman (Eds.), *Advances in neural information processing systems 19*
678 (p. 1329-1336). MIT Press.
- 679 Stramaglia, S., Cortes, J., & Marinazzo, D. (2014). Synergy and redundancy in the granger causal analysis of dynamical
680 networks. *New Journal of Physics*, 16, 105003.
- 681 Supèr, H., & Roelfsema, P. (2005). Chronic multiunit recordings in behaving animals: advantages and limitations. *Prog*
682 *Brain Res*, 147, 263-282.
- 683 Tauste Campo, A., Martinez-Garcia, M., Nàcher, V., Luna, R., Romo, R., & Deco, G. (2015). Task-driven intra-and
684 interarea communications in primate cerebral cortex. *Proc Natl Acad Sci U.S.A.*, 112(15), 4761-4766.

- 685 Thiele, A., Delicato, L., Roberts, M., & Gieselmann, M. (2006). A novel electrode-pipette design for simultaneous
686 recording of extracellular spikes and iontophoretic drug application in awake behaving monkeys. *J Neurosci Methods*,
687 *158*, 207-211.
- 688 Vinck, M., Huurdeman, L., Bosman, C., Fries, P., Battaglia, F., Pennartz, C., & Tiesinga, P. (2015). How to detect the
689 granger-causal flow direction in the presence of additive noise? *Neuroimage*, *108*, 301-18. doi:
690 10.1016/j.neuroimage.2014.12.017
- 691 Wilke, C., Ding, L., & He, B. (2008). Estimation of time-varying connectivity patterns through the use of an adaptive
692 directed transfer function. *IEEE Trans Biomed Eng*, *55*(11), 2557-2564.
- 693 Winkler, A., Ridgway, G., Webster, M., Smith, S., & Nichols, T. (2014). Permutation inference for the general linear model.
694 *Neuroimage*, *92*, 381-97. doi: 10.1016/j.neuroimage.2014.01.060
- 695 www.scipy.org. (n.d.). the scientific python library.
- 696 Zalesky, A., Fornito, A., & Bullmore, E. (2010). Network-based statistic: identifying differences in brain networks.
697 *Neuroimage*, *53*(4), 1197-1207.