

1                   

# Dependence Discovery from Multimodal Data via 2                   Multiscale Graph Correlation

3                   Cencheng Shen<sup>1</sup>, Carey E. Priebe<sup>2</sup>, Mauro Maggioni<sup>3</sup>, and Joshua T. Vogelstein\*<sup>4</sup>

4                   <sup>1</sup>Department of Statistics, Temple University

5                   <sup>2</sup>Department of Applied Mathematics and Statistics, Johns Hopkins University

6                   <sup>2</sup>Department of Mathematics, Duke University

7                   <sup>4</sup>Department of Biomedical Engineering and Institute for Computational Medicine, Johns  
8                   Hopkins University

9                   June 23, 2016

10                  

## Abstract

11                  Understanding and discovering dependence between multiple properties or measurements  
12                  is a fundamental task not just in science, but also policy, commerce, and other domains. An  
13                  ideal test for dependence would have the following properties: (1) Theoretical consistency such  
14                  that the testing power converges to 1 under any dependency structure and dimensionality. (2)  
15                  Strong empirical performance on a wide variety of low- and high-dimensional simulations. (3)  
16                  Provides insight into the nature of the dependence, rather than merely a valid p-value. (4)  
17                  On real data, detects dependence when it exists, and does not detect dependence when it  
18                  does not exist. No existing test satisfies all of these properties. In this paper we propose a  
19                  novel dependence test statistic called “Multiscale Graph Correlation” (Mgc), by combining the  
20                  ideas of distance correlation with nearest-neighbor testing. More specifically, we only use the  
21                  distance correlations amongst the nearest-neighbors of each data point, yielding a sparse,  
22                  and therefore regularized, matrix from which we can compute the test statistic. We demon-  
23                  strate that Mgc has all of the above properties via a series of theoretical proofs, numerical  
24                  simulations, and real data experiments. Specifically, we applied Mgc in several real applica-  
25                  tions: (i) detect dependence between brain disorder and hippocampus shape, (ii) determine  
26                  whether either of two pipelines can detect dependence between brain activity and personality,  
27                  and (iii) do not inflate non-existent dependence between resting activity and a spurious stim-  
28                  ulation. Mgc performs as well or better than previously proposed methods in essentially all  
29                  theory, low-dimensional and high-dimensional simulations, and real data experiments. Mgc is  
30                  therefore poised to be useful in a wide variety of applications, requiring only data and a dis-  
31                  similarity function for both measurement types. Both MATLAB and R code are provided here:  
32                  <https://github.com/jovo/RankdCorr/>.

33                  *Keywords:* testing independence, distance correlation, k-nearest-neighbor, local correlation coef-  
34                  ficient, permutation test

---

\*jovo@jhu.edu

## 35 **Contents**

36	<b>A Simulation Functions</b>	<b>19</b>
37	<b>B Dependence Measures</b>	<b>24</b>
38	B.1 (Global) MANTEL Test . . . . .	26
39	B.2 (Global) Distance Correlation . . . . .	27
40	B.3 (Global) Modified Distance Correlation . . . . .	28
41	B.4 Multiscale Graph Correlations (Mgc) . . . . .	29
42	B.5 Heller, Heller & Gorfine (HHG) . . . . .	29
43	<b>C Mgc Algorithms and Testing Procedures</b>	<b>30</b>
44	C.1 Algorithms . . . . .	31
45	C.2 Discussions of Optimal Scale Estimation . . . . .	32
46	<b>D Proofs</b>	<b>38</b>

47 Detecting dependency among multiple data sets is one of the most important and fundamental  
48 tasks in computational statistics and data science. Indeed, prior to embarking on a predictive  
49 machine learning investigation, one might first check whether any dependence is detectable; if not,  
50 high-quality predictions will be unlikely. The founders of statistics first highlighted the importance  
51 of this task, starting with Pearson, who developed Pearson's Product-Moment Correlation statistic  
52 [1]. Since then, researchers have consistently developed new and improved methods (see [2] for  
53 a recent review and discussion).

54 In the era of big data, several challenges emerge as particularly prevalent and therefore, problem-  
55 atic. First, the dependencies between different modalities of data can be highly **non-linear**. While  
56 this has always been the case, the relative abundance of data has led to an increased demand in  
57 checking for dependence in many previously uninvestigated settings. Second, the **dimensionality**  
58 of individual samples is growing at exponential rates, with genomics and connectomics data, for  
59 example, often accruing millions or billions of dimensions per data point. At the same time, the  
60 **sample sizes** are not increasing proportionally, meaning that we often have datasets with very  
61 high-dimensions and relatively low sample size. Third, the data are often **complicated**: networks,  
62 shapes, questionnaires, semi-structured text are all typical examples. For example, we may de-  
63 sire to understand whether brain shape and disease status are related, so that we can develop  
64 prognostic biomarkers to combat the deleterious effects of degenerative neurological disorders [3].  
65 Fourth, because we will often have a data deluge, with myriad different measurements, it is impor-  
66 tant to be able to compute the results reasonably **efficiently**. Fifth, when working with big data,  
67 statistical procedures often have hyper-parameters that require tuning. Many such procedures  
68 lack any guidance in choosing the value of those hyper-parameters, thereby requiring users of the  
69 procedures to concoct their own heuristics. It is desirable that a procedure is **adaptive**, in that it  
70 can automatically set its hyper-parameters in a valid way. Finally, as alluded to above, checking for  
71 dependence is rarely the final step in the analysis. Frequently, investigators and analysts desire  
72 more than a simple p-value, rather, they desire some insight into the nature of the **dependence**  
73 **structure**, which can then inform them in terms of how to proceed. We desire tests that satisfy  
74 the above desiderata, both in theory as well as in extensive simulations and real data problems.

75 There are two key insights from the literature that we combine to develop our methodology that  
76 satisfies the above desiderata. First, a collection of pairwise comparisons suffices to characterize a  
77 joint distribution [4]. Second, nonlinear manifolds can be approximated by local linear spaces. Our  
78 approach, Multiscale Graph Correlation (Mgc), leverages and improves upon recent developments

79 from both subdisciplines of data science.

80 Interpoint pairwise comparison matrices have been used for over 100 years for various statistical  
81 purposes. Perhaps one of the earliest examples comes from Karl Pearson, who created Pearson's  
82 Product-Moment Correlation [1], a special case of something subsequently named "generalized  
83 correlation coefficients" [5]. Specifically, let  $a_{ij}$  to be *some* function of  $x_i$  and  $x_j$  (for example, their  
84 Euclidean distance), and define  $b_{ij}$  similarly for pairs  $y_i$  and  $y_j$ . Thus,  $A = \{a_{ij}\}$  and  $B = \{b_{ij}\}$  are  
85 the interpoint comparison matrices for  $x$  and  $y$ , respectively. Without loss of generality, assuming  
86  $A$  and  $B$  have zero mean, the generalized correlation coefficient can then be written:

$$C = \frac{1}{z} \sum_{i,j=1}^n a_{ij} b_{ij}, \quad (1)$$

87 where  $z$  is proportional to standard deviations of  $A$  and  $B$ , that is  $z = n^2 \sigma_a \sigma_b$ . In words,  $C$   
88 is the correlation across *pairwise comparisons*, rather than the individual data samples.  $C$  has  
89 many well known special cases historically, including Pearson's correlation [1], Spearman's and  
90 Kendall's rank correlation [5], and Mantel's correlation [6]. Recently, Szekely et al. [7] extended  
91 these approaches, rather than merely subtracting the overall mean, they subtract the row means  
92 and column means, resulting in "doubly centered" distances. Impressively, they proved that this  
93 "distance correlation" (`Dcorr`) statistic is a consistent test for independence for any joint distribution  
94 (under suitable regularity conditions), that is, the `Dcorr`'s power approaches 1 as sample size  
95 approaches infinity, for any joint distribution of finite dimension and finite second moments. By  
96 adjusting the high-dimensional bias of `Dcorr`, Szekely et al. [8] further proposed `Mcorr`, which is  
97 proven to be consistent as dimensions increase to infinity as well. Existing generalized correlation  
98 coefficient based tests therefore work well in high dimensions and low sample sizes [7, 8], including  
99 in complicated domains [9], and are reasonably computationally efficient. But, they struggle in  
100 various non-linear settings, do not automatically adapt to the data, and do not provide much insight  
101 into the nature of the dependence, other than a p-value.

102 A deep insight that the generalized correlation coefficient tests have failed to capitalized on, al-  
103 though it has reaped benefits in myriad data science problems, is that of locality. Locality has been  
104 utilized for classification and regression [10], data compression [11], and recommender systems  
105 [12], to name a few. Moreover, it has become an invaluable tool in unfolding nonlinear geometry  
106 in many recent development of nonlinear embedding algorithms, dating back to the 1950s [13],  
107 and more recently making a resurgence with the advent of Isomap [14, 15], Local Linear Em-  
108 bedding [16, 17], and Laplacien eigenmaps [18], among many others. The concept of locality,  
109 while popular within certain fields has only entered into testing very infrequently [19–21]. These

110 approaches have the advantage of naturally operating on complicated data, including categorical  
 111 and structured data, as well as strong theoretical guarantees. However, these approaches focus  
 112 on two-sample testing, rather than dependence testing. Moreover, these tests all fail to provide a  
 113 convenient or automatic way to choose the appropriate neighborhood size. This omission burdens  
 114 the user of the test to come up with a heuristic, thereby greatly impairs their practical usages and  
 115 finite sample performances. In fact, adaptively choosing the scale to regularize the data is a perni-  
 116 cious problem in essentially all existing local procedures, and a good choice of joint neighborhood  
 117 can better unfold the nonlinearity and match data sets of multiple modality [22]. A primary contri-  
 118 bution of our work is an efficient mechanism both for looking across scales, and for choosing the  
 119 optimal scale.

## 120 Multiscale Graph Correlation

121 All dependence tests start from the same setting: we observe  $n$  pairs of observations  $(x_i, y_i)$ , and  
 122 we desire to know whether the  $x$ 's and  $y$ 's are independent of one another, and if so, what is the  
 123 nature of that dependence structure (Figure 1 provides an example where  $x$  and  $y$  are nonlinearly  
 124 dependent).

125 Our MGC combines generalized correlation coefficients with graph distances, in an effort to effi-  
 126 ciently uncover local relationships and optimize the independence test. Specifically, let  $R(a_{ij})$  be  
 127 the “rank” of  $x_i$  relative to  $x_j$ , that is,  $R(a_{ij}) = k$  if  $x_i$  is the  $k^{th}$  closest point (or “neighbor”) to  
 128  $x_j$ , starting from 1 to  $n$ , and define  $R(b_{ij})$  equivalently for the  $y$ 's. For any neighborhood size  $k$   
 129 around each  $x$  and any neighborhood size  $l$  around each  $y$ , we define the rank-truncated pairwise  
 130 comparisons:

$$a_{ij}^k = \begin{cases} a_{ij} - \bar{a}^k, & \text{if } R(a_{ij}) \leq k, \\ 0, & \text{otherwise;} \end{cases} \quad b_{ij}^l = \begin{cases} b_{ij} - \bar{b}^l, & \text{if } R(b_{ij}) \leq l, \\ 0, & \text{otherwise;} \end{cases} \quad (2)$$

131 where  $\bar{a}^k$  and  $\bar{b}^l$  are two mean-adjusting scalars such that  $\sum_{i,j=1}^n a_{ij}^k = \sum_{i,j=1}^n b_{ij}^l = 0$ . Then  
 132 we can define a *local* variant of any global generalized correlation coefficient, by excluding large  
 133 distances:

$$C^{kl} = \frac{1}{z_{kl}} \sum_{i,j=1}^n a_{ij}^k b_{ij}^l, \quad (3)$$

134 where  $z_{kl} = n^2 \sigma_a^k \sigma_b^l$ , with  $\sigma_a^k$  and  $\sigma_b^l$  being the standard deviations for the truncated pairwise com-  
 135 parisons. There are a maximum of  $n^2$  different local correlations, one for each possible combina-

136 tions of  $k$  and  $l$  (more technical details of Mgc are in Appendix B.4). Among all  $n^2$  local statistics,  
137  $\{C^{kl}\}$ , Mgc selects the best local statistic for testing.

138 Having defined how to compute Mgc, we face three challenges to make the method practical. First,  
139 in addition to the test statistic, we need to compute the null distribution, so that we may find the  
140 critical values and p-values. Second, naïvely, computing all local  $C^{kl}$  statistics would require an  
141 unacceptably large computational budget. Third, having computed all local statistics, we require a  
142 method for choosing the optimal neighborhood size, in such a way that the test is still consistent,  
143 and not biased (meaning that the resultant p-value is valid).

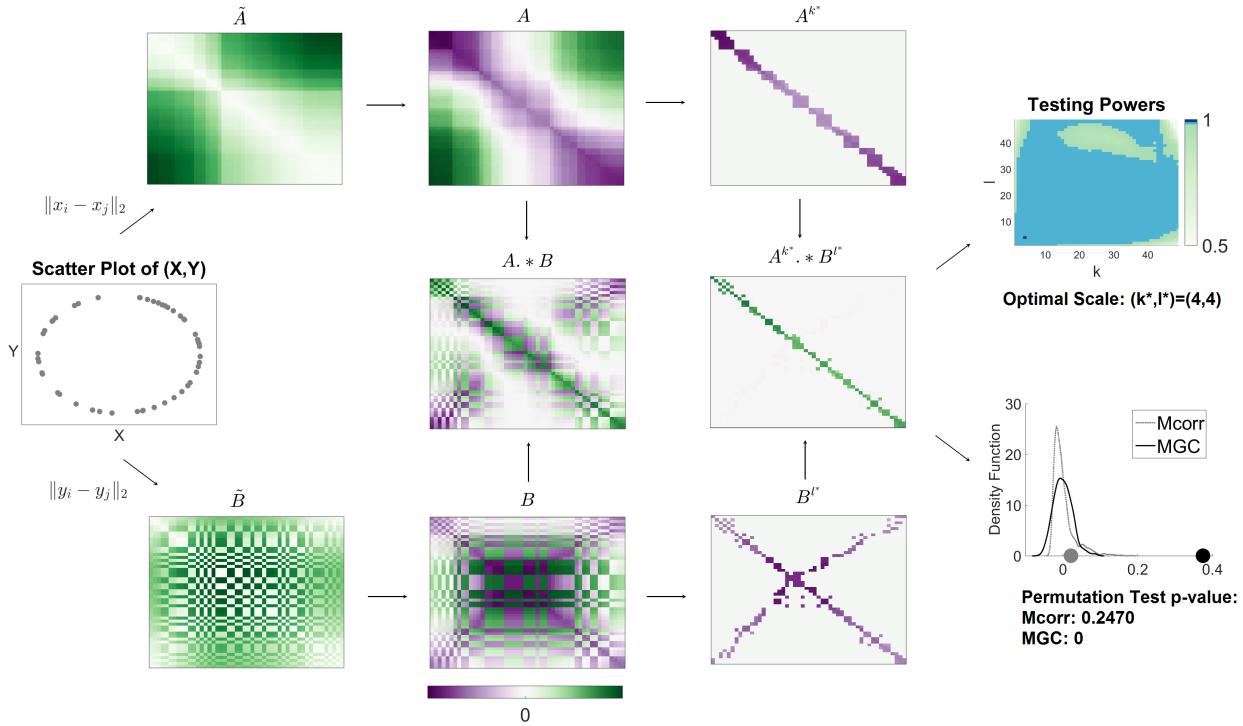
144 Computing the p-values from the test statistic is actually straightforward, thanks to the advent of  
145 permutation testing [23]. Specifically, we can permute the labels of either the  $x_i$ 's or the  $y_i$ 's, and  
146 then compute the Mgc statistics on the permuted data. By permuting the labels, we have rendered  
147 the two different views of the data independent. Doing so many times yields an empirical estimate  
148 of the null distribution, which we can use to compute the critical value and p-value. This procedure  
149 is somewhat time consuming, which makes computing the test statistics for all neighborhoods  
150 efficiently even more important.

151 Nearly all algorithms that employ *some* kind of regularization face a similar dilemma: how to  
152 efficiently choose the hyper-parameters.

153 Most manifold learning algorithms require that the user essentially runs the entire algorithm again  
154 from scratch for each different hyper-parameter setting, a pursuit that can be exponentially taxing  
155 as the number of hyper-parameters increases. In our case, once the rank information is provided,  
156 each distance-based local correlation takes  $O(n^2)$  to compute (Pseudocode 1 in Appendix C.1),  
157 which means a straightforward algorithm to compute all local correlations would take  $O(n^4)$ .

158 However, we have devised an algorithm for exactly computing *all* local correlations in  $\mathcal{O}(n^2 \log n)$ ,  
159 essentially the same running time complexity as the global correlation coefficient (the additional  
160 log factor is for sorting to find the neighbors, see Pseudocode 2 in Appendix C.1 for details). We do  
161 so by noting that the sufficient statistics for larger neighborhood sizes include those for the smaller  
162 sizes, so we can simply keep track of them as we iteratively increase neighborhood size. The end  
163 result is Mgc can be computed in comparable time as the other leading dependence tests.

164 Therefore, we can efficiently compute all local correlations for a given pair of data and the permuted  
165 data, which yields the p-values for each neighborhood size (see Pseudocode 3 in Appendix C.1  
166 for details). But it does not tell us which neighborhood sizes are optimal.



**Figure 1:** Flowchart for  $\text{Mgc}$  computation: Column 1:  $(X, Y)$  have a circle relationship. Column 2: The heat maps of  $\tilde{A}$  and  $\tilde{B}$ , which are the pairwise Euclidean distance matrices of  $X$  and  $Y$ . All distance entries are non-negative. Column 3: The top and bottom panels are the heat maps of  $A = \{a_{ij}\}$  and  $B = \{b_{ij}\}$ , which are the properly centered distance matrices of  $\tilde{A}$  and  $\tilde{B}$ . The center panel is the heatmap of the entry-wise products of  $A$  and  $B$ , summing over which yields the un-normalized  $\text{Mcorr}$  statistic. As the entries of  $A$  and  $B$  can be either positive or negative, the entry-wise products can be either positive or negative for nonlinear dependencies, which causes  $\text{Mcorr}(X, Y)$  to be close to 0 and the p-value to be in-significant, as shown in column 5. Column 4: The top and bottom panels are the heat maps of local  $A$  and  $B$ , i.e.,  $A^{k^*} = \{a_{ij}^{k^*}\}$  and  $B^{l^*} = \{b_{ij}^{l^*}\}$ , where  $(k^*, l^*) = (4, 4)$  is the optimal scale for the circle relationship. The center panel is the heatmap of the entry-wise products of local  $A$  and  $B$ , summing over which yields the un-normalized  $\text{Mgc}$  statistic  $C^*$ .  $\text{Mgc}$  successfully identifies the optimal local structure for correlation testing, and the resulting entry-wise products are dominantly non-negative, which causes  $\text{Mgc}(X, Y)$  to be much larger than 0 and the p-value to be significant, as shown in column 5. Column 5: The top panel is the testing powers of all local correlations, where the optimal scale is shown as a dark blue point with many adjacent scales being very close to optimal (light blue points). The bottom panel shows  $\text{Mgc}(X, Y)$  and  $\text{Mccorr}(X, Y)$  as dark and gray dots on the x-axis, as well as the distribution of the permuted test statistics.

167 Our procedure for estimating the optimal scale searches for regions of neighborhood sizes for  
168 which p-values are consistently low, guarding against noisy scales that appear optimal, and com-  
169 bating bias added by looking at many different scales. We assert that the optimal scale is the  
170 largest neighborhood size in that region. We define the p-value of `MGC` to be the p-value from the  
171 optimal scale, and declare significant dependency when the p-value is less than  $\alpha$ , often 0.05 (see  
172 Pseudocode 4 and 5 in Appendix C.1 for details).

173 **Theoretical Consistency of `MGC`**

The formal testing scenario is as follows: we observe  $n$  pairs of observations,  $(x_i, y_i)$ , and we desire to know whether the  $x$ 's are independent of the  $y$ 's. To cast this problem as a statistical inference query requires specifying a statistical model, that is, a collection of possible distributions from which we may assume the data arise. To make the investigation as general as possible, we consider the largest possible set of distributions: any possible joint distribution  $f_{xy}$ . If  $x$  and  $y$  were independent, then it would follow that  $f_{xy} = f_x f_y$ ; in other words, for independent data, the joint distribution is equal to the product of the marginals. Therefore, we have the following hypothesis testing scenario:

$$H_0 : f_{xy} = f_x f_y,$$
$$H_A : f_{xy} \neq f_x f_y.$$

174 The power of a test is defined as the probability that it correctly rejects the null when the null is  
175 indeed false. As defined above, a test is consistent if its power converges to 1 as sample size  
176 increases. Let  $C_t$  denote a global generalized correlation coefficient based test, that is,  $t$  might  
177 indicate Pearson, `MANTEL`, `DCORR`, or `MCORR`, and let  $\beta(C_t^*)$  denote the power of the corresponding  
178 multiscale version. Recalling from the work Szekeley et al. that `DCORR` and `MCORR` are both consis-  
179 tent tests. More specifically, `DCORR` is consistent whenever  $f_{xy}$  has finite dimension and bounded  
180 variance, and `MCORR` is consistent even as dimension increases to infinity. Denote the set of distri-  
181 butions satisfying consistency for a given test by  $\mathcal{F}_t$ , where  $t$  indicates which test we are referring  
182 to. Then, we have the following theorem:

183 **Theorem 1.**  $\beta(C_t^*) \rightarrow 1$  for all  $f_{xy}$  in  $\mathcal{F}_t$ .

184 Therefore, `MGC` is consistent against all dependent alternatives for which its global counterpart is.  
185 However, asymptotic consistency does not convey to us how quickly `MGC` achieves optimal power

186 in various settings, and whether it exhibits significant advantage over its global counterpart and  
187 other popular methods. For that, we turn to numerical simulations.

## 188 Finite Sample Simulation Experiments

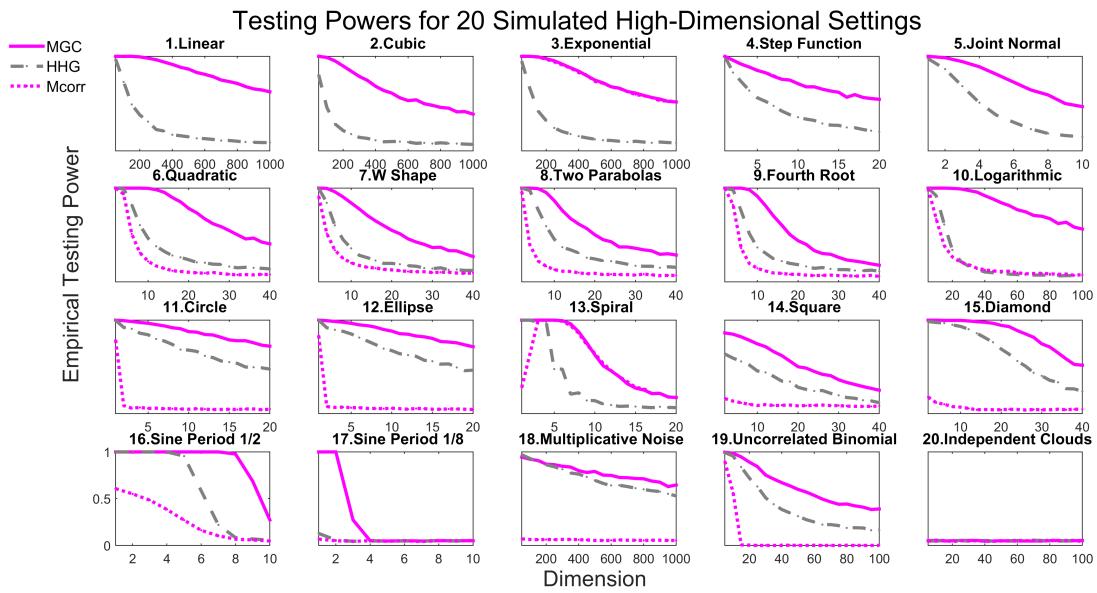
189 From this section onwards, unless mentioned otherwise, our  $M_{GC}$  is always implemented for  $MCORR$ ,  
190 due to its theoretical consistency and numerical advantages throughout.

191 Based on the previous section, we understand  $M_{GC}$  can be a consistent tests in a wide variety  
192 of settings (all finite dimensional joint distributions with bounded variance) But, our theoretical  
193 results do not shed light on the finite sample performance of  $M_{GC}$ . Specifically, we are interested  
194 in comparing our newly proposed local tests in a comprehensive set of simulations, to previous  
195 tests like  $H_{HG}$ ,  $MCORR$ ,  $DCORR$ , and  $MANTEL$ , each of which performs well in a fraction but not all of  
196 the simulations.

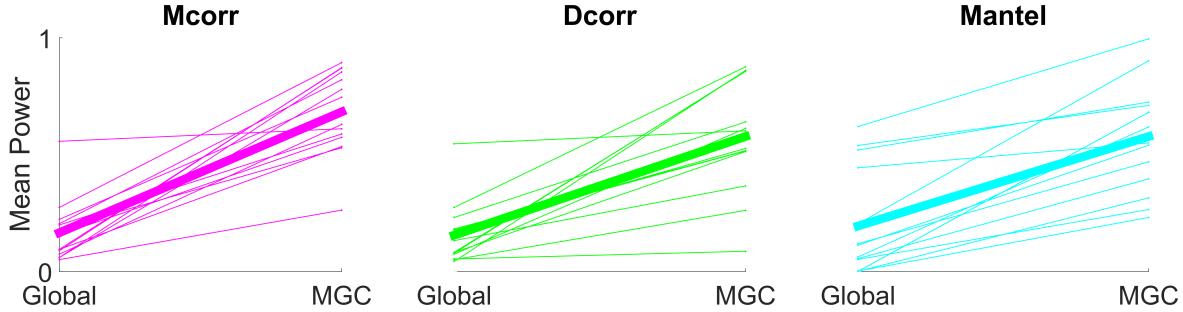
197 To do so, we consider 20 different joint distributions  $f_{xy}$ . A large fraction of these are taken exactly  
198 from existing literature [7, 24–26], and we have added several additional settings. They include:  
199 linear and nearly linear (1-5), polynomial (6-12), trigonometric (13-17), uncorrelated but nonlinearly  
200 dependent (18-19), and an independent relationship (20). Details for each setting are given in  
201 Appendix A, with a visualization of each dependency shown in Supplementary Figure A1.

202 Figure 2 shows the testing powers of  $M_{GC}$ ,  $MCORR$ , and  $H_{HG}$  versus the dimensionality of  $x$  (see  
203 Methods for details), with the sample sizes fixed at  $n = 100$  for each simulation. Note that the  
204 dimensionality of  $y$  increases in only a subset of the settings.

205 The advantage of  $M_{GC}$  over its global counterpart  $MCORR$  and  $H_{HG}$  is stark. For the nearly linear  
206 settings,  $M_{GC}$  and  $MCORR$  are essentially identical and significantly better than  $H_{HG}$  as the dimension  
207 increases. For the remaining nonlinear dependencies,  $M_{GC}$  achieves superior power than  $H_{HG}$  and  
208  $MCORR$  for all functions, often by a significant margin. For the independent simulation, all tests yield  
209 powers at the significance level  $\alpha$ , indicating no more false positives than expected according to  
210 the theory.



**Figure 2:** Powers of different methods for 20 different dependence structures, estimated by the empirical distributions of the test statistics under the null and the alternative on the basis of 10,000 Monte-Carlo replicates. 2,000 additional MC replicates are used for optimal scale estimation for Mgc. Each panel shows empirical testing power on the abscissa at a significant level  $\alpha = 0.05$ , and the dimensionality on the ordinate. Mgc empirically achieves similar or better power than the previous state of the art approaches for all sample sizes on all problems.



**Figure 3:** Average powers slopegraphs comparing global and MGC tests. For each global test, the left side corresponds to the mean power of each simulation in Figure 2, the right side corresponds to the respective MGC mean power. The thin solid lines are shown for 6-19, because MGC equals the global correlation for 1-5 and 20. Then the thick solid line summarizes how the overall mean power (including 6-19) changes from global to MGC. It is clear that MGC always significantly improves over its global counterpart.

### 211 MGC dominates Global Counterparts

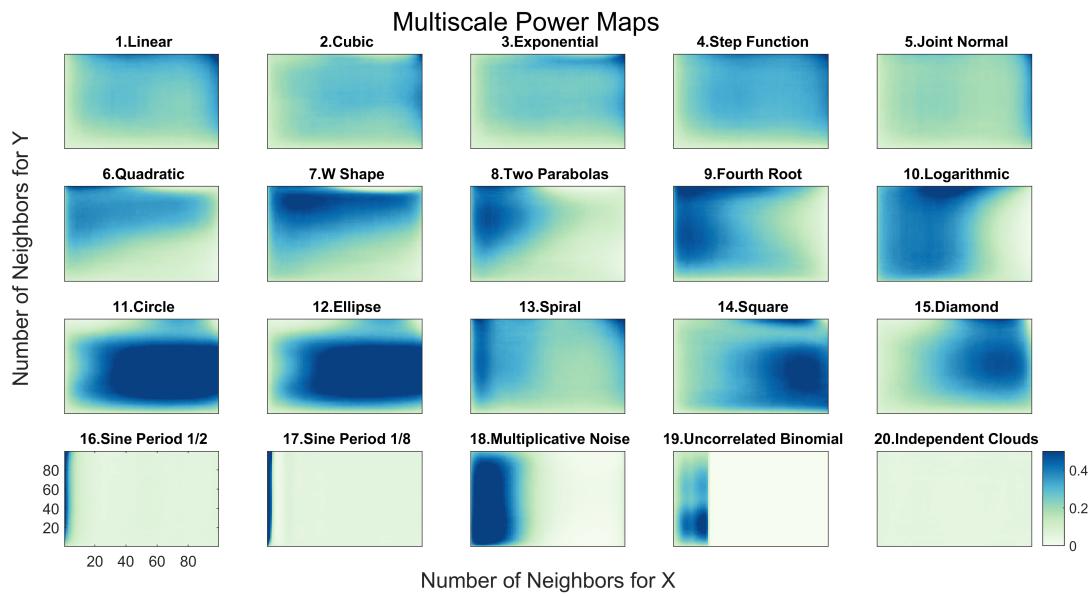
212 To better summarize the advantages of global versus local, Figure 3 shows how the powers change  
 213 from each global correlation to its MGC implementation, for each dependency in Figure 2. Indeed,  
 214 MGC always improves over its global correlation, regardless of the global correlation that is being  
 215 used. Note that the actual powers for DCORR, MANTEL, and their MGC variants are included in  
 216 Appendix A, where the same conclusion for MGC superiority still hold. We also present in Appendix  
 217 A an additional simulation setting with increasing sample size and fixed dimensionality to observe  
 218 that the powers of MGC converge to 1 faster than all the benchmarks for nearly all dependencies.

### 219 Discovery of Dependency Across Scales

220 A multiscale power map is a heatmap of powers for all neighborhood sizes, for a given joint distribu-  
 221 tion and sample size. Figure 4 provides the multiscale power maps for all 20 different scenarios for  
 222 a specified dimensionality (see caption for details), illustrating how the powers of local correlations  
 223 change with respect to increasing neighborhood sizes.

224 The multiscale power map sheds light into the intrinsic dependency structure. For nearly linear  
 225 dependencies (1-5), the best neighborhood choice is always the largest scale, i.e.,  $k = l = n$ .  
 226 For all strongly nonlinear dependencies (6-19), MGC almost always chooses a smaller scale in a  
 227 distribution dependent fashion. Furthermore, similar dependencies have similar local correlation

228 structure, and thus similar optimal scales. For example, quadratic (6) and W (7) are both polyno-  
 229 mials of degree 2 with different coefficients, and their power maps are quite similar to each other.  
 230 Similarly, (16) and (17) are the same trigonometry function (sine) with different periods, and they  
 231 share a narrow range of significant local correlations. Both circle (11) and eclipse (12), as well  
 232 as square (14) and diamond (15), are closely related functions, and have similar multiscale power  
 233 maps. Note that for almost all simulations, there exist a large portion of adjacent local neighbor-  
 234 hoods that are equally significant, which is an important observation that we use to approximate  
 235 the optimal Mgc scale for real data.



**Figure 4:** Influence of neighborhood size on testing power of local correlations at  $\alpha = 0.05$ . For each of the 20 panels, the abscissa denotes the number of neighbors for  $X$  (the scale increases from left to right), and the ordinate denotes the number of neighbors for  $Y$  (the scale increases from bottom to top). For each simulation, the sample size is  $n = 100$ , and the dimension is determined by the largest dimension for Mgc to have powers exceeding the threshold 0.5. Each different simulation yields a different surface, highlighting the importance of understanding local scale in terms of understanding the data.

236 The above described qualitative descriptions led us to believe the following two conjectures. First,  
 237 for linear dependencies, the optimal Mgc scale is the global one. Second, under certain nonlinear  
 238 dependencies, Mgc can achieve a better finite-sample testing power than its corresponding global  
 239 correlation. Indeed, we were able to prove both of these claims:

240 **Theorem 2.** *If  $x$  is linearly dependent on  $y$ , then for any  $n$  it always holds that*

$$\beta(C^{mn}) = \beta(C^*) = \beta(C). \quad (4)$$

241 Thus the optimal scale for Mgc is the global scale for linearly dependent data.

242 On the other hand, for finite sample nonlinear dependencies (which better characterize all real  
243 data) we have the following theorem.

244 **Theorem 3.** There exists  $f_{xy}$  and  $n$  such that

$$\beta(C^*) > \beta(C). \quad (5)$$

245 Thus multiscale graph correlation can be better than its global correlation coefficient under certain  
246 nonlinear dependency, for finite sample.

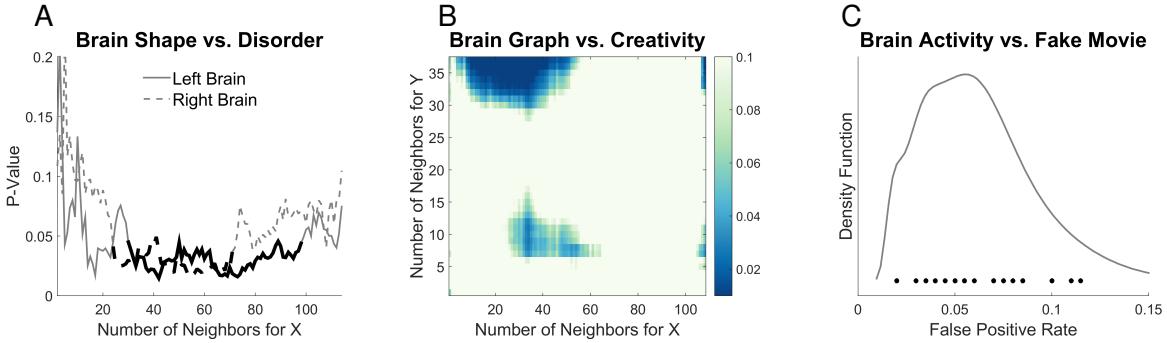
247 Note that Theorem 2 and Theorem 3 hold for any of Mgc varieties, including Dcorr, Mcorr, and  
248 Mantel. The proofs of Theorem 2 and 3 are both in Appendix D. The proof of Theorem 2 is  
249 straightforward. The proof of Theorem 3 is a constructive one. More specifically, we constructed  
250 quadratic function and sampled data a finite number of times and exactly compute the power for  
251 both Mgc and Dcorr, proving that Mgc has higher power in this setting. This shows that Mgc can  
252 outperform its global counterpart even for the most modest nonlinear functions. Because any  
253 function can be approximated by a polynomial expansion [27], the proof of Theorem 3 suggests  
254 that Mgc is able to outperform its corresponding global correlation on a wide variety of nonlinear  
255 functions, which is indeed the case throughout the numerical simulations.

## 256 Real Data Experiments

### 257 Only Local Scales can Detect Dependence

258 Our first real data experiment investigates whether brain shape and disease status are indeed  
259 dependent on one another. Previous investigations have linked major depressive disorder to the  
260 hippocampus shape [3, 28], though global tests were unable to detect a statistically significant  
261 dependence structure at the  $\alpha = 0.05$  level.

262 This brain shape versus disease dataset consists of  $n = 114$  subjects, for each we have an  
263 MRI scans as well as a categorical variable indicating whether the subject is clinically depressed,  
264 high-risk, or non-affected. From the MRI data, previous work we extracted both the left and right  
265 hippocampi. For the brain shape “view” of the data, we compute the interpoint comparison matri-  
266 ces using a nonlinear landmark matching approach [3, 29]. For the categorical disorder variable,



**Figure 5:** (A) Local correlation p-value curves with respect to  $k = 2, \dots, 114$  at  $l = 4$  for brain vs disease. Dark lines correspond to the largest region of significant scales. (B) Local correlation p-value heat map with respect to  $k = 2, \dots, 109$  and  $l = 2, \dots, 38$  for brain MIGRAINE vs CCI. (C) Density estimate for the false positive rates of Mgc on the brain vs noise experiments, with the actual rate of each data shown as dots above the x-axis.

267 we use squared Euclidean distance, then add 1 to every non-diagonal entry (so only the diagonals  
268 are of distance 0).

269 We consider two dependence tests, one for each hemisphere: is hippocampus shape independent  
270 of depressive state. Figure 5A provides the p-value curves for Mgc for  $k = 2, \dots, n$  at  $l = 4$  (we  
271 only show  $l = 4$  because the other curves look similar). Many local scales yield significant p-  
272 values (around 0.01) for both hemispheres, whereas the global scale does not detect a significant  
273 dependence in either hemisphere. None of the previously proposed dependence tests under  
274 consideration (MANTEL, DCORR, MCORR, or HHG) were able to detect dependence for both (not shown).

#### 275 MGC can provide insight into the nature of the dependence structure

276 The next real data experiment investigates whether brain networks and personalities are indepen-  
277 dent of one another. Previous work [30] investigated whether individual voxels were related to  
278 specific dimensions of personality, but were unable to compare entire brain networks to a higher-  
279 dimensional characterization of personality. Figure 5B shows that global dependence tests can  
280 ascertain whether the whole brain-network is independent of the five-factor personality traits [31].  
281 However, the global test is quite fragile, even ignoring a single subject from the global test can  
282 render the test non-significant. On the other hand, Mgc is more robust, there is a whole region  
283 of neighborhood sizes such that the test is quite significant. Moreover, that the local tests per-  
284 forms optimally with approximately 30 neighbors suggests that these data have multiple cohorts,

285 for which the dependence structure likely differs. This result therefore suggests the next investi-  
286 gatory steps to take to further understand the nature of the dependence structure between brain  
287 networks and personality.

288 **MGC Does Not Inflate False Positive Rates**

289 In the last experiment, Mgc is applied to test independence between brain voxel activities and  
290 non-existent stimulus similar to [32], by using 26 resting state fMRI data sets from the 1000 func-  
291 tional connectomes project ([http://fcon\\_1000.projects.nitrc.org/](http://fcon_1000.projects.nitrc.org/)). We used CPAC [33] to  
292 estimate regional time-series, in particular, using the sequence of pre-processing decisions de-  
293 termined to optimize discriminability [34]. The output for each dataset is the resting state fMRI  
294 time-series data containing 200 regions of interest for 200 time-steps. We then also generate an  
295 independent stimulus by sampling from a standard normal at each time step. Of course, the brain  
296 activity data and the stimuli are independent by construction. For each brain region, we test: is  
297 activity of that brain region independent of the time-varying stimuli. We pool brain activity over all  
298 of the samples from the population. Any regions that are detected significant are false positives by  
299 definition. By testing reach brain region separately, we obtain a distribution of false positive rates.  
300 If our test is unbiased, that distribution should be centered around the critical level, which we set  
301 at 0.05 for this experiment.

302 To conduct this test, we must construct a distance matrix for brain region activity, and another for  
303 the stimulus. For each brain region, we compute  $a_{ij} = \|\mathbf{x}_{\cdot i} - \mathbf{x}_{\cdot j}\|_2^2$ , for all  $(i, j)$  pairs, where  $\mathbf{x}_{\cdot i}$   
304 denotes the observation vector of all subjects at time-step  $i$ . For the stimulus, we similarly compute  
305 the Euclidean distance between activity at all pairs of time-steps:  $b_{ij} = \|y_i - y_j\|_2^2$ . Note that the  
306 distance matrices at different brain regions are distinct, but the stimulus is the same for all brain  
307 regions during the same experiment.

308 For each data set, the above test is carried out for each brain region, and the false positive rates of  
309 Mgc for each dataset are shown in Figure 5C. Mgc false positive rate is centered around the critical  
310 level 0.05, as it should be. In contrast, standard methods for fMRI analysis, such as generalized  
311 linear models, significantly increase or decrease the false discovery rates, depending on the data  
312 [32, 35].

313 **Discussion**

314 We propose multiscale graph correlation to test independence between measurement types. We  
315 demonstrate via simulations that Mgc empirically performs well in linear and non-linear settings,  
316 regardless of the dimension, sample size, and noise. Moreover, it efficiently adapts to the data, to  
317 provide not just a valid p-value, but also a picture of which scales contain the dependence struc-  
318 ture. We then prove that it achieves optimal power asymptotically no matter what the dependence  
319 structure is, even in complicated settings. In real data experiments it revealed dependence where  
320 global methods failed, revealed the locality of dependence where global methods succeeded, and  
321 did not falsely detect signals when there were none.

322 A method closely related to distance correlation tests arises from the machine learning commu-  
323 nity: kernel-based independence test [36–38]. Recent work has demonstrated the equivalence  
324 between these kernel tests and the energy statistics work [39, 40]. Thus, we may be able to glean  
325 further insights by casting Mgc within the kernel framework. Two other tests merit particular men-  
326 tion at this point. First, Dumcke et al [41] recently proposed a related nearest-neighbor based  
327 test. Unfortunately, their proposed test requires estimating relative high-dimensional densities,  
328 and therefore, does not perform particularly well, nor does it have strong theoretical support. Fi-  
329 nally, Reshef et al [42] is another competing methodology, but does not perform as well as energy  
330 based tests in various benchmarks [24], and their actual test is an approximation with unknown  
331 error bound relative to their theoretical claims.

332 Although our definition of local correlation coefficient is fast to implement, generally applicable  
333 to any global correlation, and achieves good testing powers, there are multiple ways to combine  
334 neighborhood information into a particular global correlation coefficient. So it is possible that  
335 the testing performance may be further improved, by tailoring a different centering or ranking  
336 scheme for a given global correlation, or by coming up with a different rank-truncated pairwise  
337 comparison. Overall, a more thorough investigation on the finite-sample performance of Mgc, its  
338 possible extensions, and other existing methods, are much needed in the future to enhance our  
339 understanding of dependence discovery.

340 Furthermore, the optimal scale for Mgc is also of interest, such as how to more accurately select the  
341 local scale under unknown models for a particular inference task, and the implication of the optimal  
342 scale on the geometry of underlying dependency, etc. Another direction we are investigating is  
343 how to choose the optimal metric for given data. Beyond the dependence testing framework, it may

<sup>344</sup> also be promising to pursue the applications of MGC and local correlations in other closely-related  
<sup>345</sup> subjects, such as dimension reduction, classification, other testing and prediction domains, etc.

## <sup>346</sup> References

- <sup>347</sup> [1] K. Pearson, *Proceedings of the Royal Society of London* **58**, 240 (1895). <sup>3, 4</sup>
- <sup>348</sup> [2] M. Reimherr, D. Nicolae, *Statistical Science* **28**, 116 (2013). <sup>3</sup>
- <sup>349</sup> [3] Y. Park, C. Priebe, M. Miller, N. Mohan, K. Botteron, *Journal of Biomedicine and Biotechnol-*  
<sup>350</sup> *ogy* p. 694297 (2008). <sup>3, 13</sup>
- <sup>351</sup> [4] J. Maa, D. Pearl, R. Bartoszynski, *Annals of Statistics* **24**, 1069 (1996). <sup>3</sup>
- <sup>352</sup> [5] M. G. Kendall, *Rank Correlation Methods* (London: Griffin, 1970). <sup>4</sup>
- <sup>353</sup> [6] N. Mantel, *Cancer Research* **27**, 209 (1967). <sup>4, 26</sup>
- <sup>354</sup> [7] G. Szekely, M. Rizzo, N. Bakirov, *Annals of Statistics* **35**, 2769 (2007). <sup>4, 9, 19, 27, 38</sup>
- <sup>355</sup> [8] G. Szekely, M. Rizzo, *Journal of Multivariate Analysis* **117**, 193 (2013). <sup>4, 28, 38</sup>
- <sup>356</sup> [9] R. Lyons, *Annals of Probability* **41**, 3284 (2013). <sup>4, 27</sup>
- <sup>357</sup> [10] C. Stone, *Annals of Statistics* **4**, 595 (1977). <sup>4</sup>
- <sup>358</sup> [11] I. Daubechies, *Ten lectures on wavelets* (SIAM, 1992). <sup>4</sup>
- <sup>359</sup> [12] B. Sarwar, G. Karypis, J. Konstan, J. Riedl, *ACM WebKDD 2000 Workshop* (2000). <sup>4</sup>
- <sup>360</sup> [13] W. Torgerson, *Multidimensional Scaling: I. Theory and method* (Psychometrika, 1952). <sup>4</sup>
- <sup>361</sup> [14] J. B. Tenenbaum, V. de Silva, J. C. Langford, *Science* **290**, 2319 (2000). <sup>4</sup>
- <sup>362</sup> [15] V. de Silva, J. B. Tenenbaum, *Advances in Neural Information Processing Systems* **15**, 721  
<sup>363</sup> (2003). <sup>4</sup>
- <sup>364</sup> [16] L. K. Saul, S. T. Roweis, *Science* **290**, 2323 (2000). <sup>4</sup>
- <sup>365</sup> [17] S. T. Roweis, L. K. Saul, *Journal of Machine Learning Research* **4**, 119 (2003). <sup>4</sup>
- <sup>366</sup> [18] M. Belkin, P. Niyogi, *Neural Computation* **15**, 1373 (2003). <sup>4</sup>

- 367 [19] D. Barton, F. David, *Research Papers in Statistics*, Wiley, New York (1966). 5
- 368 [20] J. Friedman, L. Rafsky, *Annals of Statistics* **11**, 377 (1983).
- 369 [21] M. Schilling, *Journal of the American Statistical Association* **81**, 799 (1986). 5
- 370 [22] C. Shen, J. T. Vogelstein, C. Priebe, *Submitted* (2016). 5
- 371 [23] P. Good, *Permutation, Parametric, and Bootstrap Tests of Hypotheses* (Springer, 2005). 6
- 372 [24] N. Simon, R. Tibshirani, available at <http://arxiv.org/abs/1401.7645> (2012). 9, 16, 19
- 373 [25] M. Gorfine, R. Heller, Y. Heller, available at <http://ie.technion.ac.il/~gorfinm/files/science6.pdf> (2012). 19, 29
- 375 [26] R. Heller, Y. Heller, M. Gorfine, *Biometrika* **100**, 503 (2013). 9, 29
- 376 [27] W. Rudin, *Real and Complex Analysis* (McGraw-Hill Education, 1986), third edn. 13
- 377 [28] J. Posener, et al., *American Journal of Psychiatry* **160**, 83 (2003). 13
- 378 [29] M. Beg, M. Miller, A. Trouv, L. Younes, *International journal of computer vision* **61**, 139 (2005).  
379 13
- 380 [30] J. Adelstein, et al., *PLoS ONE* **6**, e27633 (2011). 14
- 381 [31] R. R. Costa, & McCrae, *Neo PI-R professional manual*, vol. 396 (1992). 14
- 382 [32] A. Eklund, M. Andersson, C. Josephson, M. Johannesson, H. Knutsson, *NeuroImage* **61**, 565  
383 (2012). 15
- 384 [33] C. Craddock, et al., *Frontiers in Neuroinformatics* **42** (2015). 15
- 385 [34] S. Wang, C. E. Priebe, M. Maggioni, J. T. Vogelstein, *in preparation* (2016). 15
- 386 [35] A. Eklund, T. Nichols, H. Knutsson, *arXiv* (2015). 15
- 387 [36] A. Gretton, R. Herbrich, A. Smola, O. Bousquet, B. Scholkopf, *Journal of Machine Learning  
388 Research* **6**, 2075 (2005). 16
- 389 [37] A. Gretton, L. Gyorfi, *Journal of Machine Learning Research* **11**, 1391 (2010).
- 390 [38] A. Gretton, K. Borgwardt, M. Rasch, B. Scholkopf, A. Smola, *Journal of Machine Learning  
391 Research* **13**, 723 (2012). 16

- 392 [39] D. Sejdinovic, B. Sriperumbudur, A. Gretton, K. Fukumizu, *Annals of Statistics* **41**, 2263  
393 (2013). [16](#)
- 394 [40] A. Ramdas, S. J. Reddi, B. Pcos, A. Singh, L. Wasserman, *29th AAAI Conference on Artifi-*  
395 *cial Intelligence* (2015). [16](#)
- 396 [41] S. Dmcke, U. Mansmann, A. Tresch, *PLOS ONE* **9**, e107955 (2014). [16](#)
- 397 [42] D. Reshef, *et al.*, *Science* **334**, 1518 (2011). [16](#)

## 398 Acknowledgment

399 This work was partially supported by National Security Science and Engineering Faculty Fellow-  
400 ship (NSSEFF), Johns Hopkins University Human Language Technology Center of Excellence  
401 (JHU HLT COE), Defense Advanced Research Projects Agency's (DARPA) SIMPLEX program  
402 through SPAWAR contract N66001-15-C-4041, and the XDATA program of the Defense Advanced  
403 Research Projects Agency (DARPA) administered through Air Force Research Laboratory con-  
404 tract FA8750-12-2-0303. The authors thank Dr. Brett Mensh of Optimize Science for acting as our  
405 intellectual consigliere.

## 406 A Simulation Functions

407 We list the distributions of the 20 dependencies used in the simulations, which are based on a  
408 combination of the simulations used in [7, 24, 24, 25] but with some changes (such as the inclusion  
409 of additional noise and an extra weight vector) to better compare all methods throughout different  
410 dimensions and sample sizes.

411 For each sample  $x \in \mathbb{R}^{d_x}$ , we denote  $x^d, d = 1, \dots, d_x$  as the  $d$ th dimension of  $x$ . For the purpose  
412 of high-dimensional simulations,  $w \in \mathbb{R}^{d_x}$  is a decaying vector with  $w^d = 1/d$  for each  $d$ , such  
413 that  $w^T x$  is a 1-dimensional weighted summation of all dimensions of  $x$ , which equals  $x$  if  $d_x = 1$ .  
414 Furthermore,  $\mathcal{U}$  denotes the uniform distribution,  $\mathcal{B}$  denotes the Bernoulli distribution,  $\mathcal{N}$  denotes  
415 the normal distribution,  $u$  and  $v$  represent realizations from some auxiliary random variables,  $c$  is  
416 a scalar constant to control the noise level (which equals 1 for 1-dimensional simulations and 0

<sup>417</sup> otherwise), and  $\epsilon$  is sampled from an independent standard normal distribution unless mentioned  
<sup>418</sup> otherwise.

<sup>419</sup> For all of the below equations,  $(\mathbf{x}, \mathbf{y}) \stackrel{iid}{\sim} f_{xy} = f_{y|x}f_x$ . For each setting, we provide the space of  
<sup>420</sup>  $(\mathbf{x}, \mathbf{y})$ , and define each of the above distributions, and any additional auxiliary distributions.

1. Linear  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}$ ,

$$\mathbf{x} \sim \mathcal{U}(-1, 1)^{d_x},$$

$$\mathbf{y} = w^\top \mathbf{x} + c\epsilon.$$

2. Cubic  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}$ :

$$\mathbf{x} \sim \mathcal{U}(-1, 1)^{d_x},$$

$$\mathbf{y} = 128(w^\top \mathbf{x} - \frac{1}{3})^3 + 48(w^\top \mathbf{x} - \frac{1}{3})^2 - 12(w^\top \mathbf{x} - \frac{1}{3}) + 80c\epsilon.$$

3. Exponential  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}$ :

$$\mathbf{x} \sim \mathcal{U}(0, 3)^{d_x},$$

$$\mathbf{y} = \exp(w^\top \mathbf{x}) + 10c\epsilon.$$

4. Step Function  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}$ :

$$\mathbf{x} \sim \mathcal{U}(-1, 1)^{d_x},$$

$$\mathbf{y} = I(w^\top \mathbf{x} > 0) + \epsilon,$$

<sup>421</sup> where  $I$  is the indicator function.

5. Joint normal  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_x}$ : Let  $\rho = 1/2d_x$ ,  $I_{d_x}$  be the identity matrix of size  $d_x \times d_x$ ,  
 $J_{d_x}$  be the matrix of ones of size  $d_x \times d_x$ , and  $\Sigma = \begin{bmatrix} I_{d_x} & \rho J_{d_x} \\ \rho J_{d_x} & I_{d_x} \end{bmatrix}$ . Then let  $(u, v) \sim \mathcal{N}(0, \Sigma)$ ,  
 $\epsilon \sim \mathcal{N}(0, I_{d_x})$ ,

$$\mathbf{x} = u,$$

$$\mathbf{y} = v + 0.5c\epsilon.$$

6. Quadratic  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}$ :

$$\mathbf{x} \sim \mathcal{U}(-1, 1)^{d_x},$$

$$\mathbf{y} = (w^\top \mathbf{x})^2 + 0.5c\epsilon.$$

7. W Shape  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}$ :  $u \sim \mathcal{U}(-1, 1)^{d_x}$ ,

$$\begin{aligned}\mathbf{x} &\sim \mathcal{U}(-1, 1)^{d_x}, \\ \mathbf{y} &= 4 \left[ \left( (w^\top \mathbf{x})^2 - \frac{1}{2} \right)^2 + w^\top u / 500 \right] + 0.5c\epsilon.\end{aligned}$$

8. Two Parabolas  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}$ :  $\epsilon \sim \mathcal{U}(0, 1)$ ,  $u \sim \mathcal{B}(0.5)$ ,

$$\begin{aligned}\mathbf{x} &\sim \mathcal{U}(-1, 1)^{d_x}, \\ \mathbf{y} &= ((w^\top \mathbf{x})^2 + 2c\epsilon) \cdot (u - \frac{1}{2}).\end{aligned}$$

9. Fourth Root  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}$ :

$$\begin{aligned}\mathbf{x} &\sim \mathcal{U}(-1, 1)^{d_x}, \\ \mathbf{y} &= |w^\top \mathbf{x}|^{\frac{1}{4}} + \frac{c}{4}\epsilon.\end{aligned}$$

10. Logarithmic  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_x}$ :  $\epsilon \sim \mathcal{N}(0, I_{d_x})$

$$\begin{aligned}\mathbf{x} &\sim \mathcal{N}(0, I_{d_x}), \\ \mathbf{y}^d &= 2\log(\mathbf{x}^d) + 3c\epsilon^d,\end{aligned}$$

422 **for**  $d = 1, \dots, d_x$ .

11. Circle  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}$ :  $u \sim \mathcal{U}(-1, 1)^{d_x}$ ,  $\epsilon \sim \mathcal{N}(0, I_{d_x})$ ,  $r = 1$ ,

$$\begin{aligned}\mathbf{x}^d &= r \left( \sin(\pi u^{d+1}) \prod_{j=1}^d \cos(\pi u^j) + 0.4\epsilon^d \right) \text{ for } d = 1, \dots, d_x - 1, \\ \mathbf{x}^{d_x} &= r \left( \prod_{j=1}^{d_x} \cos(\pi u^j) + 0.4\epsilon^{d_x} \right), \\ \mathbf{y} &= \sin(\pi u^1).\end{aligned}$$

423 12. Ellipse  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}$ : Same as above except  $r = 5$ .

13. Spiral  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}$ :  $u \sim \mathcal{U}(0, 5)$ ,  $\epsilon \sim \mathcal{N}(0, 1)$ ,

$$\begin{aligned}\mathbf{x}^d &= u \sin(\pi u) [\cos(\pi u)]^d \text{ for } d = 1, \dots, d_x - 1, \\ \mathbf{x}^{d_x} &= u [\cos(\pi u)]^{d_x}, \\ \mathbf{y} &= u \sin(\pi u) + 0.4(d_x - 1)\epsilon.\end{aligned}$$

14. Square  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_x}$ : Let  $u \sim \mathcal{U}(-1, 1)$ ,  $v \sim \mathcal{U}(-1, 1)$ ,  $\epsilon \sim \mathcal{N}(0, 1)^{d_x}$ ,  $\theta = -\frac{\pi}{8}$ . Then

$$\begin{aligned}\mathbf{x}^d &= u \cos \theta + v \sin \theta + 0.05d_x \epsilon^d, \\ \mathbf{y}^d &= -u \sin \theta + v \cos \theta,\end{aligned}$$

424 **for**  $d = 1, \dots, d_x$ .

425 15. Diamond  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_x}$ : Same as above except  $\theta = -\frac{\pi}{4}$ .

16. Sine Period 1/2  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}$ :  $u \sim \mathcal{U}(-1, 1)$ ,  $v \sim \mathcal{N}(0, 1)^{d_x}$ ,  $\theta = 4\pi$ ,

$$\mathbf{x}^d = u + 0.02d_x v^d \text{ for } d = 1, \dots, d_x,$$

$$\mathbf{y} = \sin(\theta x) + c\epsilon.$$

426 17. Sine Period 1/8  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}$ : Same as above except  $\theta = 16\pi$  and the noise is changed  
427 to  $0.5c\epsilon$ .

18. Multiplicative Noise  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_x}$ :  $u \sim \mathcal{N}(0, I_{d_x})$ ,  $\epsilon \sim \mathcal{N}(0, I_{d_x})$ ,

$$\mathbf{x} \sim \mathcal{N}(0, I_{d_x}),$$

$$\mathbf{y}^d = u^d \mathbf{x}^d + 0.5\epsilon^d,$$

428 for  $d = 1, \dots, d_x$ .

19. Uncorrelated Binomial  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}$ :  $u \sim \mathcal{B}(0.5)$ ,

$$\mathbf{x} \sim \mathcal{B}(0.5)^{d_x},$$

$$\mathbf{y} = (2u - 1)w^\top \mathbf{x} + 0.6\epsilon.$$

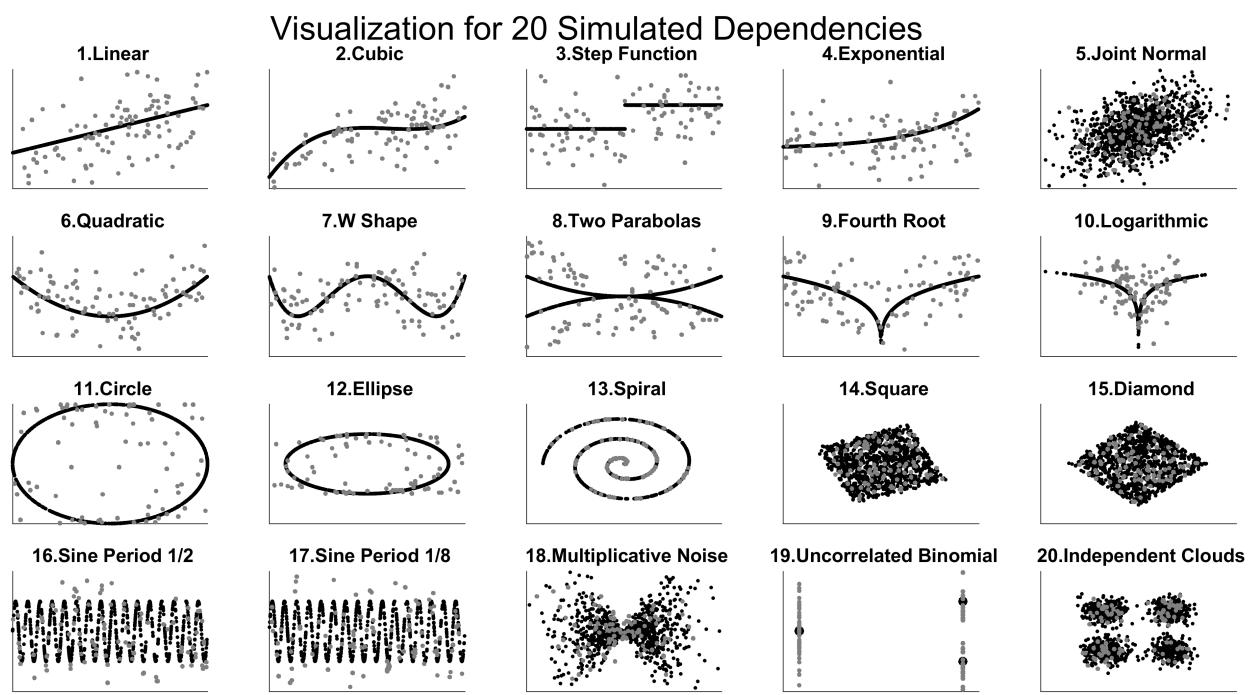
20. Independent Clouds  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_x}$ : Let  $u \sim \mathcal{N}(0, I_{d_x})$ ,  $v \sim \mathcal{N}(0, I_{d_x})$ ,  $u' \sim \mathcal{B}(0.5)^{d_x}$ ,  
 $v' \sim \mathcal{B}(0.5)^{d_x}$ . Then

$$\mathbf{x} = u/3 + 2u' - 1,$$

$$\mathbf{y} = v/3 + 2v' - 1.$$

429 For each distribution,  $x$  and  $y$  are clearly dependent except (20); for some settings (11-15) they  
430 are conditionally independent upon conditioning on the respective auxiliary variables, while for  
431 others they are "directly" dependent. Then we can independently generate  $(x_i, y_i)$  from  $(\mathbf{x}, \mathbf{y})$  for  
432  $i = 1, \dots, n$ , set  $X = [x_1, \dots, x_n] \in \mathbb{R}^{d_x \times n}$  and  $Y = [y_1, \dots, y_n] \in \mathbb{R}^{d_y \times n}$ , and calculate local /  
433 global correlations for the sample data. A visualization of each dependency is shown in Figure A1.

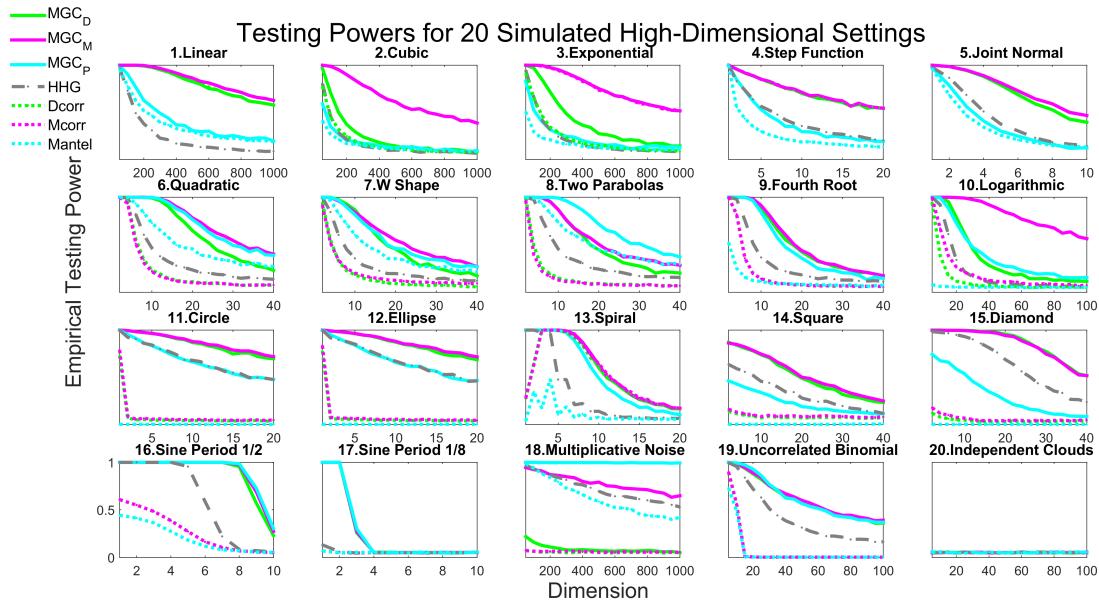
434 For the increasing dimension simulation in the main paper, we always set  $c = 0$  and  $n = 100$ ,  
435 with  $d_x$  increasing while  $d_y = d_x$  for type 5, 10, 14, 15, 18, 20 and  $d_y = 1$  otherwise. The decaying  
436 vector  $w$  is utilized for  $d_x > 1$  to treat higher dimensions as small perturbations, which creates a  
437 meaningful setting for testing power comparison. The powers of all three Mgc implementations in  
438 this setting are provided in Figure A2, where we denote  $Mgc_D$  as the Mgc for Dcorr,  $Mgc_M$  as the  
439 Mgc for Mcorr,  $Mgc_P$  as the Mgc for Mantel.



**Figure A1:** Visualization of the 20 dependencies for 1-dimensional simulations. The blue points are generated with noise ( $c=1$ ) at  $n = 100$  to show the actual sample data in testing, and the red points are generated without noise at  $n = 1000$  to highlight each underlying dependency.

440 Here we also present an additional setting, which sets  $d_x = d_y = 1$  and  $c = 1$  with the sample size  
 441  $n$  increasing from 5 to 100. The parameter before  $c$  (e.g., there is a 80 before  $c$  in type 2) is a tuned  
 442 noise parameter for some dependencies, so the testing powers can be compared meaningfully  
 443 for each simulation, i.e., in the absence of noise, the testing powers may converge to 1 at very  
 444 small  $n$  for some trivial dependencies like linear; and it is also more meaningful to consider noisy  
 445 simulations in practice. The powers of all methods in this setting are provided in Figure A3, with  
 446 the multiscale power maps shown in Figure A4.

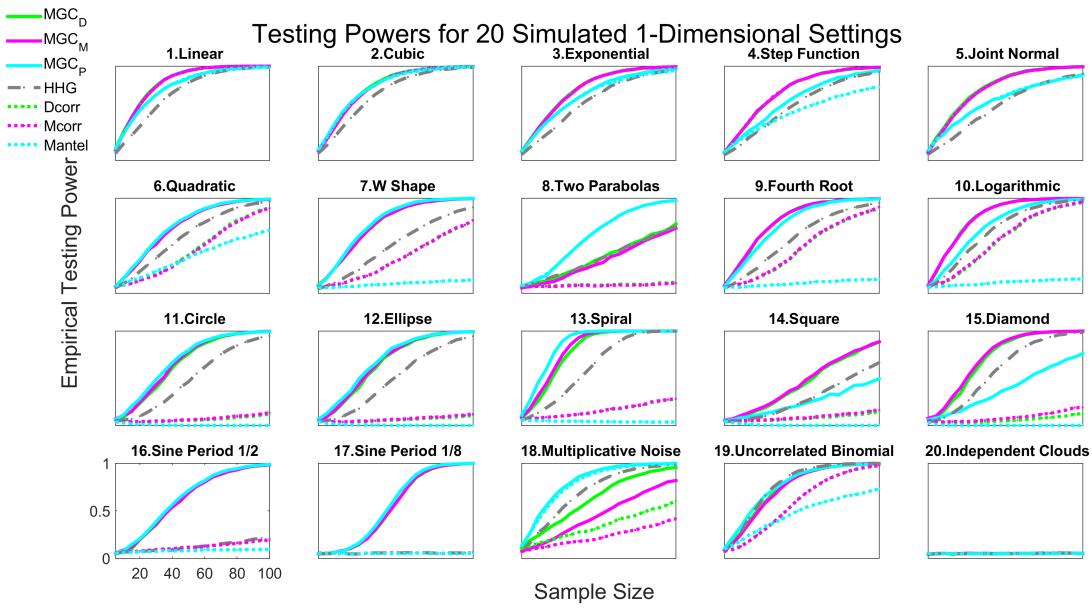
447 Clearly Mgc always improves over its global counterpart, and always has a large advantage re-  
 448 gardless of the underlying dependency structure, the dimensionality, the sample size, or noise.



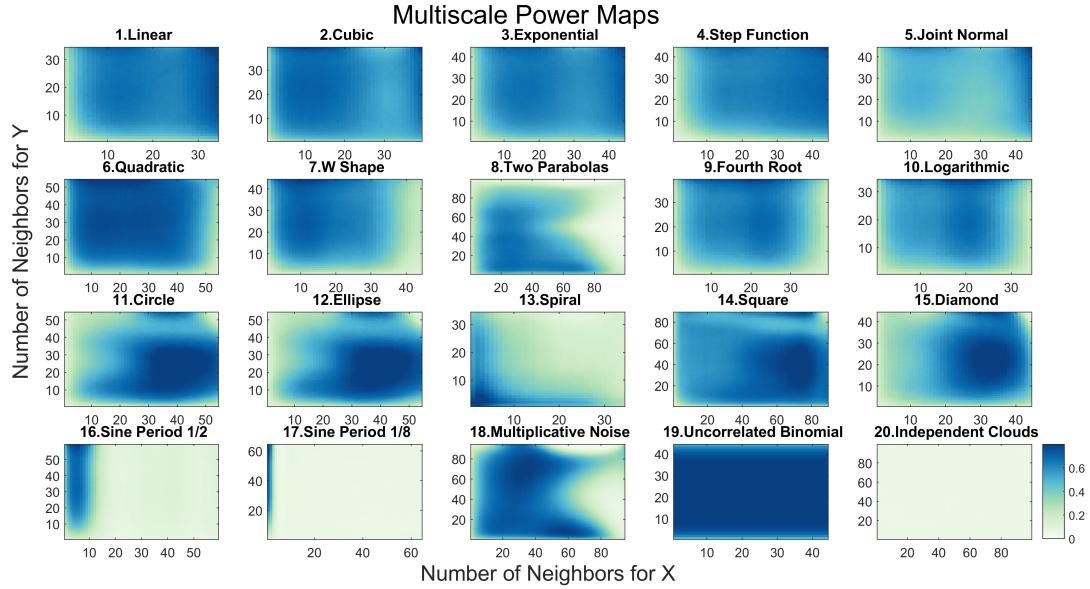
**Figure A2:** Same as Figure 2 but includes all three different Mgc implementations.

## 449 **B Dependence Measures**

450 In this section, we review the MANTEL test, distance correlation, modified distance correlation, the  
 451 Mgc statistic, and the HHG statistic in order. Note that for Dcorr / Mcorr, we implement them in a  
 452 slightly different but equivalent way from the original definition.



**Figure A3:** Powers of different methods for 20 different 1-dimensional dependence structures, estimated by the empirical distributions of the test statistics under the null and the alternative on the basis of 10,000 Monte-Carlo replicates. 2,000 additional MC replicates are used for optimal scale estimation for  $\text{Mgc}$ . Each panel shows empirical testing power on the abscissa at a significant level  $\alpha = 0.05$ , and sample size on the ordinate.  $\text{Mgc}$  empirically achieves similar or better power than the previous state of the art approaches for all sample sizes on nearly all problems.



**Figure A4:** Influence of neighborhood size on testing power of local correlations. For each simulation, the dimension is 1, and the sample size is determined by the first sample size  $n$  for MGC to have powers exceeding the threshold 0.8.

### 453 B.1 (Global) MANTEL Test

454 Given the Euclidean distance matrices  $\tilde{A}$  and  $\tilde{B}$ , the MANTEL coefficient [6] is defined as

$$\text{Mantel}(X, Y) = \frac{\sum_{i \neq j}^n (a_{ij} - \bar{a})(b_{ij} - \bar{b})}{\sqrt{\sum_{i \neq j}^n (a_{ij} - \bar{a})^2 \sum_{i \neq j}^n (b_{ij} - \bar{b})^2}}, \quad (1)$$

455 where  $A = \tilde{A}$ ,  $B = \tilde{B}$ ,  $\bar{a} = \frac{1}{n(n-1)} \sum_{i \neq j}^n (a_{ij})$  and similarly for  $\bar{b}$ . Then the MANTEL test is carried out  
456 by the permutation test.

457 Unlike distance correlation and H<sub>HG</sub>, the MANTEL test is not consistent against all dependent alter-  
458 natives, but it has been a very popular method in biology and ecology due to its simplicity. It is  
459 clear from Figure A2 and A3 that global MANTEL is sub-optimal and appears to be not consistent  
460 for many dependencies, yet MGC<sub>P</sub> achieves comparable performances as other variants of MGC,  
461 which implies that MGC<sub>P</sub> may be consistent against most, if not all dependent alternatives.

462 **B.2 (Global) Distance Correlation**

463 Given two distance matrices  $\tilde{A}$  and  $\tilde{B}$  of the sample data  $X$  and  $Y$ , the sample distance covariance  
 464 is defined by doubly centering the distance matrices:

$$dcov(X, Y) = \frac{1}{n^2} \sum_{i,j=1}^n a_{ij} b_{ij}, \quad (2)$$

where  $A = H\tilde{A}H$ ,  $B = H\tilde{B}H$  with  $H = I_n - \frac{J_n}{n}$ . Then the sample distance variance is defined as

$$\begin{aligned} dvar(X) &= \frac{1}{n^2} \sum_{i,j=1}^n a_{ij}^2, \\ dvar(Y) &= \frac{1}{n^2} \sum_{i,j=1}^n b_{ij}^2, \end{aligned}$$

465 and the sample distance correlation equals

$$Dcorr(X, Y) = \frac{dcov(X)}{\sqrt{dvar(X) \cdot dvar(Y)}}. \quad (3)$$

466 It is shown in [7] that as  $n \rightarrow \infty$ ,  $Dcorr(X, Y) \rightarrow Dcorr(\mathbf{x}, \mathbf{y}) \geq 0$ , where  $Dcorr(\mathbf{x}, \mathbf{y})$  denotes  
 467 the population distance correlation between the underlying random variable  $\mathbf{x}$  and  $\mathbf{y}$ . The pop-  
 468 ulation distance correlation is defined by the characteristic functions, which is 0 if and only if  $\mathbf{x}$   
 469 and  $\mathbf{y}$  are independent. Thus the sample distance correlation is a consistent statistic for testing  
 470 independence, i.e., the testing power  $\beta_\alpha(Dcorr(X, Y))$  converges to 1 as  $n$  increases, at any type  
 471 1 error level  $\alpha$ . Note that all of  $dcov$ ,  $dvar$ ,  $Dcorr$  are always non-negative; and the consistency  
 472 result assumes finite second moments of  $\mathbf{x}$  and  $\mathbf{y}$ , which holds for a family of metrics not limited  
 473 to the Euclidean distance [9]. Also note that the  $Dcorr$  above is actually the square of distance  
 474 correlation in [7], but for ease of presentation the square naming is dropped here.

475 Alternatively, calculating the distance covariance by  $A = H\tilde{A}$  and  $B = \tilde{B}H$  gives the same statis-  
 476 tic as in Equation 2, i.e., instead of using doubly centered distance matrices, it is the same to  
 477 singly center one distance matrix by row and the other distance matrix by column. Then  $Dcorr$  by  
 478 singly centered distance matrices has the same testing power as the original  $Dcorr$ , because dis-  
 479 tance covariance is equivalent to distance correlation in the permutation test (note that the actual  
 480  $Dcorr$  statistic by single centering is different from the original  $Dcorr$ , as using single centering  
 481 changes the distance variances).

482 In our implementation of global / local  $Dcorr$ , we always use singly centered distance matrices  
 483 rather than doubly centered distance matrices. Although they are equivalent for the testing power

484 of global  $\text{DCORR}$ , our alternative implementation improves the testing power of local  $\text{DCORR}$  and  
 485  $\text{MGC}$ . This is because the ranking information of  $\tilde{A}$  and  $\tilde{B}$  are better preserved in singly centered  
 486 distance matrices, so that  $\text{MGC}$  is more effective in excluding far-away points that exhibit insignificant  
 487 dependency. This applies to  $\text{MCORR}$  as well.

### 488 B.3 (Global) Modified Distance Correlation

489 In case of high-dimensional data where the dimension  $d_x$  or  $d_y$  increases with the sample size  $n$ ,  
 490 the sample distance correlation may no longer be appropriate. For example, even for independent  
 491 Gaussian distributions,  $\text{Dcorr}(X, Y) \rightarrow 1$  as  $d_x, d_y \rightarrow \infty$ , which may severely impair the testing  
 492 power of sample  $\text{DCORR}$  in high-dimensional simulations.

493 The modified distance correlation is proposed in [8] to tackle the bias of sample  $\text{DCORR}$ . Denote  
 494 the Euclidean distance matrices as  $\tilde{A}$  and  $\tilde{B}$ , the doubly centered distance matrices as  $\hat{A}$  and  $\hat{B}$ ,  
 495 the modified distance covariance is defined as

$$496 \quad mcov(X, Y) = \frac{n}{(n-1)^2(n-3)} \left( \sum_{i \neq j}^n a_{ij} b_{ij} - \frac{2}{n-2} \sum_{j=1}^n a_{jj} b_{jj} \right), \quad (4)$$

496 where  $A$  modifies the entries of  $\hat{A}$  by

$$a_{ij} = \begin{cases} \hat{a}_{ij} - \frac{\bar{\hat{a}}_{ij}}{n}, & \text{if } i \neq j, \\ \frac{n \sum_i \hat{a}_{ij} - \sum_{i,j} \hat{a}_{ij}}{n^2}, & \text{if } i = j, \end{cases}$$

497 and so is  $B$ . Then  $mvar(X)$  and  $mvar(Y)$  can be similarly defined.

498 If  $mvar(X) \cdot mvar(Y) \leq 0$ , the modified distance correlation is set to 0 (negativity can only occur  
 499 when  $n \leq 2$ , equality can only happen in some special cases); otherwise it is defined as

$$500 \quad \text{Mcorr}(X, Y) = \frac{mcov(X, Y)}{\sqrt{mvar(X) \cdot mvar(Y)}}. \quad (5)$$

500 It is shown in [8] that  $\text{Mcorr}(X, Y)$  is an unbiased estimator of the population distance correlation  
 501  $\text{Dcorr}(x, y)$  for all  $d_x, d_y, n$ ; and  $\text{MCORR}$  is approximately normal even if  $d_x, d_y \rightarrow \infty$ . Thus it is a  
 502 consistent statistic for testing independence, but may work better than  $\text{DCORR}$  under high-dimension  
 503 dependencies.

504 Similar to the alternative implementation of  $\text{DCORR}$ , we can also use singly centered distance ma-  
 505 trices for  $\hat{A}$  and  $\hat{B}$  in defining  $\text{MCORR}$ , which does not alter the theoretical advantages of original  
 506  $\text{MCORR}$ . We further set  $A_{ii} = B_{ii} = 0$  for all  $i$ , which simplifies the expression of  $\text{MCORR}$  and is  
 507 asymptotically equivalent for the testing purpose.

508 **B.4 Multiscale Graph Correlations ( $M_{GC}$ )**

509 For any generalized correlation coefficient, its local correlations can be directly implemented as  
510 in Equation 3, by plugging in the respective  $a_{ij}$  and  $b_{ij}$  from Equation 1 and sorting the distance  
511 matrices column-wise as in Equation 2.

512 In particular, **MANTEL** sets  $a_{ij}$  and  $b_{ij}$  as the respective entry of  $\tilde{A}$  and  $\tilde{B}$  (the Euclidean distances).  
513 **Dcorr** lets  $a_{ij}$  and  $b_{ij}$  be the respective matrix entry of  $A$  and  $B$  (the doubly centered distance  
514 matrices), then the sample means  $\bar{a}, \bar{b}$  are automatically 0. **Mcorr** slightly modifies  $a_{ij}$  and  $b_{ij}$  of  
515 **Dcorr** to adjust their high-dimensional bias. As discussed already, our version of  $M_{GC_M}$  is based  
516 on single centering throughout: we take  $a_{ij} = b_{ij} = 0$  when  $i = j$ , otherwise set  $a_{ij}$  as the matrix  
517 entry of  $H\tilde{A} - \tilde{A}/n$ , and set  $b_{ij}$  as the entry of  $\tilde{B}H - \tilde{B}/n$ . Then the local version of **Mcorr** follows  
518 by Equation 3.

519 Generally, there are a total of  $\max(R(a_{ij})) \times \max(R(b_{ij}))$  local correlations, which equals  $n^2$  when  
520 there exists no repeating data. Note that we use minimal ranks in sorting when ties occur, which  
521 indexes all local correlations more conveniently than breaking ties randomly or using average /  
522 max ranks.

523 Among all possible local correlations, MGC picks the optimal local correlation that yields the best  
524 testing power. The optimal scale clearly exists, but is distribution dependent and is almost always  
525 non-unique. Among all local correlations, it suffices to exclude  $C^{1l}$  and  $C^{k1}$  for testing and optimal  
526 scale estimation: since  $C^{1l} = C^{k1} = C^{11}$ , they do not include any neighbor other than each obser-  
527 vation itself, merely count the diagonal terms in the distance matrices, and are not meaningful for  
528 the testing purpose.

529 **B.5 Heller, Heller & Gorfine ( $H_{HG}$ )**

530 The  $H_{HG}$  statistic applies Pearson's chi-square test to ranks of distances within each column, and is  
531 shown to be better than many global tests including **Dcorr** under common nonlinear dependencies  
532 in [25, 26]. Like **Dcorr** and **Mcorr**,  $H_{HG}$  is distance-based and consistent, but not in the form of the  
533 generalized correlation coefficient; and like our  $M_{GC}$ , it makes use of the rank information, but in a  
534 distinct manner.

Given the Euclidean distance matrices  $\tilde{A} = [\tilde{a}_{ij}]$  and  $\tilde{B} = [\tilde{b}_{ij}]$ , we denote

$$\begin{aligned} H_{11}(i, j) &= \sum_{q=1, q \neq i, j}^n I(\tilde{a}_{ik} \leq \tilde{a}_{ij}) I(\tilde{b}_{ik} \leq \tilde{b}_{ij}) \\ H_{12}(i, j) &= \sum_{q=1, q \neq i, j}^n I(\tilde{a}_{ik} \leq \tilde{a}_{ij}) I(\tilde{b}_{ik} > \tilde{b}_{ij}) \\ H_{21}(i, j) &= \sum_{q=1, q \neq i, j}^n I(\tilde{a}_{ik} > \tilde{a}_{ij}) I(\tilde{b}_{ik} \leq \tilde{b}_{ij}) \\ H_{22}(i, j) &= \sum_{q=1, q \neq i, j}^n I(\tilde{a}_{ik} > \tilde{a}_{ij}) I(\tilde{b}_{ik} > \tilde{b}_{ij}), \end{aligned}$$

and the  $\text{H}_{\text{HG}}$  statistic is defined as

$$\text{H}_{\text{HG}}(X, Y) = \sum_{i=1, j \neq i}^n \frac{(n-2)(H_{12}(i, j)H_{21}(i, j) - H_{11}(i, j)H_{22}(i, j))^2}{H_{1.}(i, j)H_{2.}(i, j) - H_{.1}(i, j)H_{.2}(i, j)},$$

- 535 where  $H_{1.} = H_{11} + H_{12}$ ,  $H_{2.} = H_{21} + H_{22}$ ,  $H_{.1} = H_{11} + H_{21}$ , and  $H_{.2} = H_{12} + H_{22}$ . It is clear  
 536 that  $\text{H}_{\text{HG}}$  is structurally different from **Dcorr** / **Mcorr** / **MANTEL**, cannot be conveniently expressed by  
 537 Equation 1, and there is no direct extension of local correlation to  $\text{H}_{\text{HG}}$ .  
 538 The permutation test using the  $\text{H}_{\text{HG}}$  statistic is consistent against all dependent alternatives. In  
 539 our numerical simulations,  $\text{H}_{\text{HG}}$  falls a bit short when testing against high-dimensional and noisy  
 540 linear dependencies, but is often more advantageous than global correlations under nonlinear  
 541 dependencies, which makes it a strong competitor in general.

## 542 C Mgc Algorithms and Testing Procedures

- 543 In this section we elaborate on the algorithms for computing local correlation and **Mgc**, as well as  
 544 their testing procedures in simulations and real data experiment.  
 545 Five algorithms are presented in section C.1: given the choice of a global correlation coefficient,  
 546 algorithm 1 computes one local correlation coefficient at a given  $(k, l)$ ; then algorithm 2 shows  
 547 how to compute all local correlations simultaneously; algorithm 3 computes the p-values of all  
 548 local correlation by the random permutation test; algorithm 4 approximates the optimal scale for  
 549 **Mgc** based on the p-values of all local correlations, and outputs the approximated p-value of **Mgc**;  
 550 algorithm 5 estimates the testing powers of all local statistics based on a given joint distribution  
 551 or multiple pairs of data, which can be used to more accurately estimate the optimal scale for

552 MGC when the underlying model is known or training data are given. More detailed discussions  
553 regarding the optimal scale approximation is offered in section C.2.

554 **C.1 Algorithms**

555 All algorithms are implemented in Matlab and R with the pseudo-code shown below. For ease of  
556 presentation, we assume there are no repeating data and take D<sub>corr</sub> as the global correlation in  
557 the pseudo-code.

558 Algorithm 1 shows a straightforward computation of one local correlation coefficient, which re-  
559 quires  $O(n^2)$  once the rank information is provided. This is suitable for MGC computation when  
560 the optimal local scale is known or already estimated. But using algorithm 1 to compute all local  
561 correlations would require iterating through all possible neighborhoods  $(k, l)$ , which takes  $O(n^4)$   
562 and would make the optimal scale estimation computationally inefficient.

563 To facilitate the optimal scale estimation, algorithm 2 provides a fast method to compute all lo-  
564 cal correlations in  $O(n^2)$ . An important observation is that each product  $a_{ij}b_{ij}$  is included in  $C^{kl}$   
565 if and only if  $(k, l)$  satisfies  $k \leq R(a_{ij})$  and  $l \leq R(b_{ij})$ , so it suffices to iterate through  $a_{ij}b_{ij}$  for  
566  $i, j = 1, \dots, n$ , and add the product simultaneously to all  $C^{kl}$  whose scales are no more than  
567  $(R(a_{ij}), R(b_{ij}))$ . However, accessing and adding multiple  $C^{kl}$  at the same time is not computa-  
568 tionally efficient; instead, for each product, we only add it to  $C^{kl}$  at  $(k, l) = (R(a_{ij}), R(b_{ij}))$  (so only one  
569 local scale is accessed for each operation), iterate through all products for  $i, j = 1, \dots, n$ , then add  
570 up adjacent  $C^{kl}$  for  $k, l = 1, \dots, n$ . Thus all local correlations can be computed in  $O(n^2)$ , which  
571 has the same running time complexity as the global distance correlation. There are two additional  
572 overheads: sorting the distance matrices column-wise takes  $O(n^2 \log n)$ , and properly centering  
573 the distance matrices takes  $O(n^2)$ .

574 Algorithm 3 computes the p-values of all local correlation by the permutation test with  $r$  random  
575 permutations, which takes  $O(rn^2 \log n)$ .

576 Algorithm 4 approximates the optimal scale  $(k^*, l^*)$  from the p-values of all local correlations,  
577 and outputs the approximated MGC p-value. This is necessary for testing on one pair of data  
578 with unknown model, while algorithm 5 is more appropriate for known model. Conceptually, the  
579 algorithm first searches for a set of “valid” adjacent rows  $\mathcal{K} = \{k_1, k_1 + 1, \dots, k_2 - 1, k_2\}$  such that  
580 the median p-value of  $\{p_{kl}, k \in \mathcal{K}, l = 2, \dots, n\}$  is no larger than  $\alpha/(n-1) * |\mathcal{K}|$ , otherwise we take

581  $\mathcal{K} = \{n\}$ ; and similarly determine the set of valid columns  $\mathcal{L}$ . Once  $\mathcal{K}$  and  $\mathcal{L}$  are determined, the  
582 optimal scale  $(k^*, l^*)$  is found by the scale that minimizes the p-value within  $\{p_{kl}, k \in \mathcal{K}, l \in \mathcal{L}\}$ .  
583 Clearly if the majority p-values of all local correlations are less than  $\alpha$ , then  $\mathcal{K} = \mathcal{L} = \{1, \dots, n\}$ ,  
584 and the optimal scale equals the scale that minimizes the p-values among all local correlations;  
585 if there is no valid rows and columns, then  $\text{Mgc}$  takes the largest scale and equals the global  
586 correlation. Note that the actual algorithm is a simpler version of the above description: instead of  
587 considering all possible sets of rows and check the validity, we limit the check to the most likely set  
588 of rows, by first looking for the row scale of the smallest p-value, then including all adjacent rows  
589 whose minimal p-value on the row is no larger than  $\alpha$ ; similarly for the set of columns.

590 Algorithm 5 computes the testing powers of all local correlations by repeated simulating samples  
591 generated from the joint distribution  $f_{xy}$ . Sample data under the null and the alternative are re-  
592 peatedly generated for  $r$  Monte-Carlo replicates, and algorithm 2 is applied to compute the sample  
593 local correlations under the null and the alternative. Then the testing power at each local corre-  
594 lation can be estimated, and the  $\text{Mgc}$  optimal scale can be found by maximizing the powers. This  
595 algorithm is also applicable if there exists multiple pairs of data with unknown model but similar  
596 dependency structure, then the alternative statistic can be computed from each data pair while  
597 the null statistic can be computed from each data pair under permutation. The running time is  
598  $O(rn^2 \log n)$ .

## 599 C.2 Discussions of Optimal Scale Estimation

600 To evaluate  $\text{Mgc}$  in simulations or real data, the optimal scale for  $\text{Mgc}$  always needs to be estimated  
601 first. Algorithm 5 computes the testing powers of all local correlations for known model, so the  
602 optimal scale  $(k^*, l^*)$  can be directly estimated by maximizing the testing powers (if there are more  
603 than one optimal scales, one may pick the scale that maximizes the mean difference of the test  
604 statistic under the null and the alternative). Once the optimal scale is determined, the testing  
605 power of  $\text{Mgc}$  under the given model can be quickly determined by algorithm 5, and its p-value for  
606 testing on a particular pair of data can be determined by algorithm 3.

607 If there is only one pair of data  $(X, Y)$  with unknown distributions, we have to approximate the  
608 optimal scale by algorithm 4. It makes use of Bonferroni correction to separately verify the set of  
609 rows and columns, which guarantees the false positive rate to be no higher than  $\alpha$ ; otherwise the  
610 scale is set to the largest, which guarantees the approximated  $\text{Mgc}$  is at least as powerful as the

---

**Algorithm 1** Local Correlation Computation for One Scale

---

Input: A pair of distance matrices  $(\tilde{A}, \tilde{B}) \in \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times n}$ , and the given local scale  $(k, l) \in \mathbb{R} \times \mathbb{R}$ .  
Output: The local correlation coefficient  $C^{kl} \in [-1, 1]$  at the given  $(k, l)$ .

```
1: function LOCALCORR( $\tilde{A}, \tilde{B}, k, l$ )
2:   initialize  $C^{kl}, V_k^A, V_l^B, E_k^A, E_l^B$  as 0.
3:   for  $Z := A, B$  do  $R^Z = \text{SORT}(\tilde{Z})$  end for            $\triangleright$  column-wise sorting and assume no ties
4:   for  $Z := A, B$  do  $Z = \text{CENTER}(\tilde{Z})$  end for        $\triangleright$  proper centering of the distance matrices
5:   for  $i, j = 1, \dots, n$  do
6:      $C^{kl} = C^{kl} + A_{ij}B_{ij}\mathbf{I}(R_{ij}^A \leq k)\mathbf{I}(R_{ij}^B \leq l)$            $\triangleright$  store local distance covariance
7:      $V_k^A = V_k^A + A_{ij}^2\mathbf{I}(R_{ij}^A \leq k)$            $\triangleright$  store local distance variance for  $X$ 
8:      $V_l^B = V_l^B + B_{ij}^2\mathbf{I}(R_{ij}^B \leq l)$            $\triangleright$  store local distance variance for  $Y$ 
9:      $E_k^A = E_k^A + A_{ij}\mathbf{I}(R_{ij}^A \leq k)$            $\triangleright$  store the sample means
10:     $E_l^B = E_l^B + B_{ij}\mathbf{I}(R_{ij}^B \leq l)$ 
11:   end for
12:    $C^{kl} = (C^{kl} - E_k^A E_l^B / n^2) / \sqrt{(V_k^A - E_k^{A2} / n^2)(V_l^B - E_l^{B2} / n^2)}$        $\triangleright$  normalize the local covariances
13: end function
```

---

---

**Algorithm 2**  $O(n^2 \log n)$  Algorithm for Computing All Local Correlations

---

Input: A pair of distance matrices  $(\tilde{A}, \tilde{B}) \in \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times n}$ .  
Output: All local correlation coefficients  $C^{kl} \in [-1, 1]^{n \times n}$  for  $k, l = 1, \dots, n$ .

```

1: function LOCALCORR( $\tilde{A}, \tilde{B}$ )
2:   initialize  $C$  as a zero matrix of size  $n \times n$ ;  $V^A, V^B, E^A, E^B$  as zero vectors of size  $n$ .
3:   for  $Z := A, B$  do  $R^Z = \text{SORT}(\tilde{Z})$  end for
4:   for  $Z := A, B$  do  $Z = \text{CENTER}(\tilde{Z})$  end for
5:   for  $i, j = 1, \dots, n$  do
6:      $k = R_{ij}^A$ 
7:      $l = R_{ij}^B$ 
8:      $C^{kl} = C^{kl} + A_{ij}B_{ij}$ 
9:      $V_k^A = V_k^A + A_{ij}^2$ 
10:     $V_l^B = V_l^B + B_{ij}^2$ 
11:     $E_k^A = E_k^A + A_{ij}$ 
12:     $E_l^B = E_l^B + B_{ij}$ 
13:   end for
      ▷ the next two for loops with respect to the scales guarantee the computation of all local
      covariance / variance in  $O(n^2)$ 
14:   for  $k = 1, \dots, n - 1$  do
15:      $C^{1,k+1} = C^{1,k} + C^{1,k+1}$ 
16:      $C^{k+1,1} = C^{k+1,1} + C^{k+1,1}$ 
17:     for  $Z := A, B$  do  $V_{k+1}^Z = V_k^Z + V_{k+1}^Z$  end for
18:     for  $Z := A, B$  do  $E_{k+1}^Z = E_k^Z + E_{k+1}^Z$  end for
19:   end for
20:   for  $k, l = 1, \dots, n - 1$  do
21:      $C^{k+1,l+1} = C^{k+1,l} + C^{k,l+1} + C^{k+1,l+1} - C^{k,l}$ 
22:   end for
23:   for  $k, l = 1, \dots, n$  do                                ▷ normalize all local covariances
24:      $C^{kl} = (C^{kl} - E_k^A E_l^B / n^2) / \sqrt{(V_k^A - E_k^A)^2 / n^2 (V_l^B - E_l^B)^2 / n^2}$ 
25:   end for
26: end function

```

---

---

**Algorithm 3** P-value Computation for All Local Correlations

---

Input: A pair of distance matrices  $(\tilde{A}, \tilde{B}) \in \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times n}$ , the number of permutations  $r$ .  
Output: The p-value matrix  $P \in [0, 1]^{n \times n}$  for all local distance correlations.

```
1: function PERMUTATIONTEST( $\tilde{A}, \tilde{B}, r$ )
2:    $C^{kl} = \text{LOCALCORR}(\tilde{A}, \tilde{B})$                                  $\triangleright$  calculate the observed local correlations
3:   for  $j = 1, \dots, r$  do
4:      $\pi = \text{RANDPERM}(n)$                                           $\triangleright$  generate a random permutation of size  $n$ 
5:      $C_0^{kl}[j] = \text{LOCALCORR}(\tilde{A}, \tilde{B}(\pi, \pi))$            $\triangleright$  calculate the permuted test statistics
6:   end for
7:   for  $k, l = 1, \dots, n$  do
8:      $P_{kl} = \sum_{j=1}^r (C^{kl} < C_0^{kl}[j]) / r$                    $\triangleright$  get the p-value at each local scale
9:   end for
10:  end function
```

---

611 global correlation. Still, algorithm 4 is a heuristic approach to approximate the optimal local scale,  
612 which does not guarantee the optimal local correlation to be always correctly identified.

613 To better justify algorithm 4, we compare the estimated  $M_{GC}$  power by algorithm 4 to the true  
614  $M_{GC}$  power by algorithm 5, with the global  $MCORR$  and  $H_{HG}$  as benchmarks. For each type of depen-  
615 dency in the simulation section, we generate 1,000 pairs of dependent data by the same low- and  
616 high-dimensional settings as in Figure A3 and A2; and for each pair of data, all local p-values are  
617 calculated by 1,000 random permutations. By using the true optimal scale (from the simulation sec-  
618 tion) consistently for each data pair, the true  $M_{GC}$  p-value can be computed; by using algorithm 4 to  
619 approximate the optimal scale for each pair of data separately, the estimated  $M_{GC}$  p-value can be  
620 computed; and the p-values of global  $MCORR$  and  $H_{HG}$  can also be derived. The null is rejected when  
621 the p-value is less than 0.05, and the power equals the percentage of correct rejection. Based on  
622 the powers of true  $M_{GC}$  / estimated  $M_{GC}$  /  $MCORR$  /  $H_{HG}$  shown in Figure A5, we observe that although  
623 the estimated  $M_{GC}$  power by algorithm 4 can be lower than the true  $M_{GC}$  power, it is almost always  
624 better than global  $MCORR$  and  $H_{HG}$ , and combines the better performance of the two benchmarks.

625 Note that it is tempting to directly use the optimal scale that minimizes all local p-values without  
626 the validation by algorithm 4, or generate random samples based on the given data pair and use  
627 algorithm 5 by bootstrap. However, both approaches are biased such that the false positive rate will  
628 be higher than the type 1 error in the absence of dependency. This is because for a given pair of

---

**Algorithm 4** Optimal Local Scale Approximation by P-values

---

Input: The p-value matrix  $P \in \mathbb{R}^{n \times n}$  of all local distance correlations, the type 1 error level  $\alpha$ .

Output: The approximated MGC optimal scale  $(k^*, l^*)$ , and the approximated MGC p-value  $p$ .

```
1: function MGCSCALEVERIFY( $P, \alpha$ )
2:    $\mathcal{K} = \text{VERIFYRow}(P, \alpha)$                                  $\triangleright$  search for a set of valid row indices
3:    $\mathcal{L} = \text{VERIFYRow}(P^T, \alpha)$                                 $\triangleright$  search for a set of valid column indices
4:    $[k^*, l^*] = \arg \min_{\{k \in \mathcal{K}, l \in \mathcal{L}\}} P_{kl}$             $\triangleright$  find the optimal scale within the valid range
5:    $p = P_{k^*l^*}$ 
6: end function
```

Input: Same as MGCSCALEVERIFY.

Output: The indices of valid rows.

```
1: function VERIFYRow( $P, \alpha$ )
2:   initialize  $\mathcal{K}$  as an empty set
3:    $[k^*, l^*] = \arg \min_{k,l} \{P_{kl}, k, l = 2, \dots, n\}$ 
4:   for  $k = k^*, \dots, 2$  do                                      $\triangleright$  check all row scales no larger than  $k^*$ 
5:     if  $\min\{P_{kl}, l = 2, \dots, n\} > \alpha$  then
6:       break
7:     end if
8:      $\mathcal{K} = [k, \mathcal{K}]$ 
9:   end for
10:  for  $k = k^* + 1, \dots, m$  do                                 $\triangleright$  check all row scales larger than  $k^*$ 
11:    if  $\min\{P_{kl}, l = 2, \dots, n\} > \alpha$  then
12:      break
13:    end if
14:     $\mathcal{K} = \{\mathcal{K}, k\}$ 
15:  end for
16:  if  $\text{MEDIAN}(P_{kl}, k \in \mathcal{K}, l = 2, \dots, n) > \alpha * \frac{|\mathcal{K}|}{n-1}$  then
17:     $\mathcal{K} = \{n\}$            $\triangleright$  take the largest scale if the median p-value is not sufficiently small
18:  end if
19: end function
```

---

---

**Algorithm 5** Testing Powers Computation for All Local Correlations

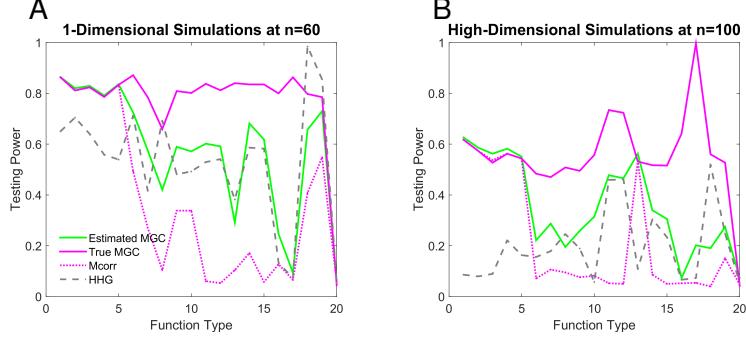
---

**Input:** A joint distribution  $f_{xy}$ , the sample size  $n$ , the number of MC replicates  $r$ , and the type 1 error level  $\alpha$ .

**Output:** The power matrix  $\beta_\alpha \in [0,1]^{n \times n}$  for all local correlations, and the Mgc optimal scale  $(k^*, l^*) \in \mathbb{R} \times \mathbb{R}$ .

```
1: function TESTINGPOWERS( $f_{xy}, n, r, \alpha$ )
2:   for  $j = 1, \dots, r$  do
3:     for  $i := [n]$  do  $(X_i^1, Y_i^1) \stackrel{iid}{\sim} f_{xy}$  end for            $\triangleright$  generate dependent samples
4:     for  $i := [n]$  do  $X_i^0 \stackrel{iid}{\sim} f_x$  end for                    $\triangleright$  generate independent samples
5:     for  $i := [n]$  do  $Y_i^0 \stackrel{iid}{\sim} f_y$  end for
6:     for  $Z := A, B$  do  $\tilde{Z}_1 = \text{DIST}(Z_1)$  end for     $\triangleright$  the distance matrices under the alternative
7:     for  $Z := A, B$  do  $\tilde{Z}_0 = \text{DIST}(Z_0)$  end for       $\triangleright$  the distance matrices under the null
8:      $C_1^{kl}[j] = \text{LOCALCORR}(\tilde{A}_1, \tilde{B}_1)$         $\triangleright$  calculate all local correlations under the alternative
9:      $C_0^{kl}[j] = \text{LOCALCORR}(\tilde{A}_0, \tilde{B}_0)$         $\triangleright$  calculate all local correlations under the null
10:   end for
11:   for  $k, l = 1, \dots, n$  do
12:      $c_\alpha = \text{CDF}_{1-\alpha}(C_{kl}^0[j], j \in [r])$            $\triangleright$  get the critical value by the empirical cumulative
        distribution under the null at each scale
13:      $\beta_\alpha^{kl} = \sum_{j=1}^r (C_{kl}^1[j] > c_\alpha) / r$             $\triangleright$  estimate the power
14:   end for
15:    $(k^*, l^*) = \arg \max(\beta_\alpha^{kl})$                        $\triangleright$  find the optimal local scale
16: end function
```

---



**Figure A5:** Comparing estimated MGC power to true MGC power, for the 1-dimensional and high-dimensional simulations. (A) 1-dimensional simulations, where  $d_x = 1$  and the sample size is chosen by the power threshold 0.8 as in Figure A4. (B) High-dimensional simulations, where  $n = 100$  and the dimension is chosen by the power threshold 0.5 as in Figure 4. The estimated MGC power by the approximated optimal scale is almost always better than global M<sub>CORR</sub> and H<sub>HG</sub>, combines the better performance of the two benchmarks, is quite close to the true MGC power, and does not inflate false signals.

629 data, a non-optimal scale can happen to have a significant p-value, which may be falsely identified  
 630 as optimal if we directly minimize all local p-values. Those erroneous scales often still exist after  
 631 a straightforward re-sampling, so random samples have the same problem. More investigations  
 632 into the bias and better methods for searching the optimal scale are two worthwhile directions for  
 633 future works.

## 634 D Proofs

635 **Theorem 1.**  $\beta(C_t^*) \rightarrow 1$  for all  $f_{xy}$  in  $\mathcal{F}_t$ .

636 *Proof.* For any  $f_{xy}$ , the power of multiscale graph correlation satisfies

$$\beta(C^*) = \max_{\mathbf{x}, \mathbf{l}} \{\beta(C^{kl})\} \geq \beta(C), \quad (6)$$

637 at any type 1 error level  $\alpha$ . So  $\beta(C^*) \rightarrow 1$  if  $\beta(C) \rightarrow 1$ .

638 Therefore  $\beta(C_t^*) \rightarrow 1$  for all  $f_{xy}$  in  $\mathcal{F}_t$ . In particular, M<sub>GCD</sub> and M<sub>GCM</sub> are consistent against all alter-  
 639 native of finite second moments, because D<sub>CORR</sub> and M<sub>CORR</sub> are consistent against all alternatives  
 640 of finite second moments by [7, 8].  $\square$

641 **Theorem 2.** If  $x$  is linearly dependent on  $y$ , then for any  $n$  it always holds that

$$\beta(C^{nn}) = \beta(C^*) = \beta(C). \quad (7)$$

642 Thus the optimal scale for MGC is the global scale for linearly dependent data.

643 *Proof.* To show that MGC is equivalent to the global correlation coefficient, it suffices to show the  
644 p-value of  $C^{kl}$  is always no less than the p-value of  $C$  for all  $k, l$  under linear dependence.

645 Under linear dependency, for any global correlation coefficient satisfying Equation 1, by Cauchy-  
646 Schwarz inequality it follows that

$$1 = C(X, Y) \geq C(X, YQ) \quad (8)$$

647 for any permutation matrix  $Q$ , where the equality holds if and only if  $X$  is a scalar multiple of  $YQ$ .

648 It follows that the p-value of  $C$  is 0, which is at the minimal.

649 Therefore the p-value of  $C^{kl}$  cannot be less than the p-value of  $C$  under linear dependency, such  
650 that the global correlation is the optimal scale for MGC under linear dependency.  $\square$

651 **Theorem 3.** There exists  $f_{xy}$  and  $n$  such that

$$\beta(C^*) > \beta(C). \quad (9)$$

652 Thus multiscale graph correlation can be better than its global correlation coefficient under certain  
653 nonlinear dependency, for finite sample.

654 *Proof.* We give a simple discrete example of  $f_{xy}$  at  $n = 7$ , such that the p-value of  $MGC_M$  is strictly  
655 lower than the p-value of MCORR.

Suppose under the alternative, each pair of observation  $(x, y)$  is sampled as follows:

$$\begin{aligned} x &\in \{-1, -\frac{2}{3}, -\frac{1}{3}, 0, \frac{1}{3}, \frac{2}{3}, 1\} \text{ without replacement,} \\ y &= x^2, \end{aligned}$$

656 which is a discrete version of the quadratic relationship in the simulations.

657 At  $n = 7$ , we can directly calculate  $C^{kl}(X, Y)$  and  $\{C^{kl}(X, YQ)\}$  for all permutation matrices  $Q$ . It  
658 follows that the p-value of MCORR is  $\frac{151}{210}$ , while  $C^{kl}(X, Y) = \frac{17}{70}$  at  $(k, l) = (2, 4)$ . Note that in this

659 case  $k$  is bounded above by  $n = 7$  while  $l$  is bounded above by 4 due the the repeating points in  
660  $Y$ .

661 Then by choosing  $\alpha = 0.25$ ,  $\text{Mgc}$  has power 1 while global  $\text{MCORR}$  has power 0, i.e.,  $\text{Mgc}$  successfully  
662 identifies the dependency in this example while global  $\text{MCORR}$  fails.

663 Note that we can always consider sample points in  $[-1, 1]$  for  $X$ , increase  $n$  and reach the same  
664 conclusion with more significant p-values; but the computation of all possible permuted test statis-  
665 tics becomes more time-consuming as  $n$  increases. The same conclusion also holds for  $\text{Mgc}_D$  and  
666  $\text{Mgc}_P$  using the same example. □