# Identifying disease-associated biomarker network features through conditional graphical model

**Shanghong Xie[1]** ⬤ | **Xiang Li[2]** | **Peter McColgan[3]** | **Rachael I. Scahill[3]** |
**Donglin Zeng[4]** ⬤ | **Yuanjia Wang[1,5]** ⬤

[1]Department of Biostatistics, Mailman School of Public Health, Columbia University, New York

[2]Statistics and Decision Sciences, Janssen Research & Development, LLC, Raritan, New Jersey

[3]Huntington's Disease Centre, Department of Neurodegenerative Disease, UCL Institute of Neurology, London, UK

[4]Department of Biostatistics, University of North Carolina, Chapel Hill, North Carolina

[5]Department of Psychiatry, Columbia University Medical Center, New York

**Correspondence**

Shanghong Xie, Department of Biostatistics, Mailman School of Public Health, Columbia University, New York, NY.
Email: sx2168@cumc.columbia.edu

## Abstract

Biomarkers are often organized into networks, in which the strengths of network connections vary across subjects depending on subject-specific covariates (eg, genetic variants). Variation of network connections, as subject-specific feature variables, has been found to predict disease clinical outcome. In this work, we develop a two-stage method to estimate biomarker networks that account for heterogeneity among subjects and evaluate network's association with disease clinical outcome. In the first stage, we propose a conditional Gaussian graphical model with mean and precision matrix depending on covariates to obtain covariate-dependent networks with connection strengths varying across subjects while assuming homogeneous network structure. In the second stage, we evaluate clinical utility of network measures (connection strengths) estimated from the first stage. The second-stage analysis provides the relative predictive power of between-region network measures on clinical impairment in the context of regional biomarkers and existing disease risk factors. We assess the performance of proposed method by extensive simulation studies and application to a Huntington's disease (HD) study to investigate the effect of HD causal gene on the rate of change in motor symptom through affecting brain subcortical and cortical gray matter atrophy connections. We show that cortical network connections and subcortical volumes, but not subcortical connections are identified to be predictive of clinical motor function deterioration. We validate these findings in an independent HD study. Lastly, highly similar patterns seen in the gray matter connections and a previous white matter connectivity study suggest a shared biological mechanism for HD and support the hypothesis that white matter loss is a direct result of neuronal loss as opposed to the loss of myelin or dysmyelination.

**KEYWORDS**

graphical model, gray matter network, Huntington's disease, mediation analysis, regularized regression

## 1 | INTRODUCTION

Biomarkers such as brain imaging measures may organize into networks of connected regions (eg, Alexander-Bloch *et al.*, 2013). The morphological features of some brain regions (eg, thickness and volume of gray matter) covary with other regions (Alexander-Bloch *et al.*, 2013). Such covariation patterns are referred to as structural covariance networks that may characterize coordinated patterns in morphology between anatomically connected regions of interest (ROIs)

(He *et al.*, 2008). Recent studies have suggested the presence of substantial heterogeneity of covariance network connections between individuals and subgroups of individuals. For example, the strengths of brain cortical region connections were observed to vary with age (Chen *et al.*, 2011) and the structural covariation patterns were altered across the lifespan (Alexander-Bloch *et al.*, 2013). In our motivating study of Huntington's disease (HD; Paulsen *et al.*, 2014), the networks of high-risk group and low-risk group are different in their edge strengths (Hamming distance is 4.4; details in Section 4).

The patterns of heterogeneous network connections among regions have been shown to predict disease phenotypes. Alexander-Bloch *et al.* (2013) suggest that behavioral and cognitive abilities are associated with the between-region covariation patterns in addition to the variability within the regions. These findings suggest that network connectivities between brain regions can be predictive of disease clinical outcomes beyond the regional imaging measures. Thus, it is desirable to evaluate the relative clinical utility of network connections in context of the regional biomarkers and other existing disease risk factors.

There is no existing method readily available to estimate heterogeneous gray matter structural covariance networks and their effects on a disease outcome. First, unlike the white matter connectivity network that can be obtained by probabilistic tractography in each individual, gray matter cortical thickness or subcortical volume of each region is measured only once from each individual and the current methods for constructing structural covariance network are at the population level, by calculating the partial correlation between brain regions over the population (Alexander-Bloch *et al.*, 2013). By leveraging between-subject variability and modeling dependence between network strength and individual covariates (eg, age, genetic variants), subject- or subgroup-specific gray matter network connections may be obtained. Second, there is a lack of methods to estimate heterogeneous networks in a multidimensional setting. When the network is assumed to be homogenous across individuals, Gaussian graphical models (Friedman *et al.*, 2008) have been extensively used to estimate a high-dimensional network defined by the inverse covariance matrix (ie, precision matrix). Fused graphical lasso (FGL; Danaher *et al.*, 2014) was proposed to jointly estimate graphical models for multiple distinct classes but assume that the classes are known. To address between-subject variability, Yin and Li (2011), Cai *et al.* (2012), and Chen *et al.* (2016) proposed to adjust the mean of variables in the network for covariates but still assume a constant precision matrix. Cheng *et al.* (2014) incorporate covariates directly into an Ising model of the network but under some modeling assumptions. Third, to the best of our knowledge, no existing method exploits heterogeneous network strength to construct subject-specific networks

and use them as potential intermediate phenotypes on disease outcomes.

A straightforward method to estimate subject-dependent network is to stratify the samples into subgroups, estimate a network by Gaussian graphical model (eg, graphical lasso; Friedman *et al.*, 2008) for each group separately, and associate the edge effects of each group with disease outcomes. However, it is unknown how the subgroups should be formed. Furthermore, such a method is inefficient when the number of subgroups is large and the number of subjects in each group is small. In our motivating study, subjects were stratified into high-risk group and low-risk group (Web Figure S1). The Hamming distance, defined as the sum of the absolute edgewise differences between two adjacency matrices, between the networks of the two groups is 4.4, and the weighted correlation is 0.87. Thus, although the connection weights of the two networks vary, they manifest correlated edge effects. These results suggest a more efficient approach to model the network edge strengths as a function of covariates to provide subject-dependent networks.

In this paper, we propose a two-stage method to estimate biomarker networks that account for heterogeneity among subjects and evaluate network's association with disease clinical outcome. In the first stage, we propose a conditional Gaussian graphical model to capture covariate-specific networks with connection strengths depending on a subject's covariates while assuming homogenous network structure, in which both the mean and precision matrix of the graphical model depend on individual covariates. To handle multidimensional parameter space, regularization is imposed to introduce model sparsity and stabilize estimation. The first-stage model will provide subject-specific network measures (edge connection strengths) between biomarkers. The identified network can be dense, but not all of the network connections are associated with clinical impairment that reduces their clinical utility. Thus, in the second stage, we use a penalized regression that includes existing disease risk factors (covariates), regional biomarkers, and between-region connections to simultaneously examine their clinical utility. The goal in the second stage is to further evaluate relative predictive power of network measures on clinical impairment compared to using regional biomarkers and covariates alone.

Our method makes several contributions. The first-stage model captures heterogeneity of network connection strength without assuming that the connections are directly measured on each individual. Such effects are captured by covariates and represent some smooth changes over the covariate space and subgroups of subjects. The second-stage model simultaneously considers a large number of biomarkers and connections as intermediate measures without first taking any transformation or dimension reduction. Third, since it is unknown which connections between biomarkers are associated with clinical impairment (eg, a lower connectivity between certain brain

regions may associate with poor cognition), our method identifies important network features with additive clinical utility in the context of existing risk factors.

We conduct extensive simulation studies to examine the performance of proposed method with varying sample sizes and number of biomarkers. We then apply the method to a large, long-term, natural history study of premanifest HD patients to investigate the effect of HD causal gene on the motor function deterioration through affecting brain subcortical and cortical gray matter atrophy networks (Paulsen *et al.*, 2014). We construct subject-dependent gray matter cortical thinning networks and subcortical volumetric networks and organize the ROIs into distinct modules. We evaluate the predictive power of the network connections on motor impairment in HD patients, and validate findings in an independent study (McColgan *et al.*, 2017). Lastly, we compare the identified HD gray matter networks with white matter connectivity networks to examine whether white matter loss is a direct result of neuronal loss (caused by the dying back of axons) or loss of myelin or dysmyelination. Our results provide new biological evidence supporting the former hypothesis. We conclude this paper with some remarks and extensions.

## 2 | METHODOLOGY

Let $X_i$ denote a $q$-dimensional vector of disease risk factor (eg, covariates including genetic variants, baseline clinical measures) of the $i$th subject. Let $M_i$ denote a vector of $p$-dimensional biomarkers. Let $Y_i$ denote a clinical outcome of interest. In our application in Section 4, $X_i$ includes Deoxyribonucleic acid (DNA) structural variation at a causal gene and baseline clinical measures, $M_i$ includes cortical and subcortical brain atrophy volumetric measures at ROIs available from structural magnetic resonance imaging (MRI), and $Y_i$ is the rate of change of motor symptoms. The structural covariation network is defined as the precision matrix of $M_i$. An overview of schematic diagram of our method is shown in Figure 1. Specifically, the first-stage model estimates subject-dependent networks of $M_i$ given $X_i$ and connections between components of $M_i$ from a conditional Gaussian graphical model. The second-stage model identifies which network connections between nodes have incremental effects on clinical outcomes in addition to $M_i$s and $X_i$. This model also estimates the effect of $X_i$ on $Y_i$ through connections between $M_i$.

### 2.1 | First-stage conditional network model

A conditional Gaussian graphical model for the distribution of $M_i$ given $X_i$, where both the mean and precision matrix of $M_i$ depend on $X_i$, can be expressed as

$$P(\boldsymbol{M}_i|\boldsymbol{X}_i) \propto \exp\left(\boldsymbol{\kappa}_i^T \boldsymbol{M}_i - \frac{1}{2}\boldsymbol{M}_i^T \boldsymbol{\Omega}_i \boldsymbol{M}_i\right), \quad (1)$$

where the $j$th element of $\boldsymbol{\kappa}_i$ is $\kappa_{ij} = \boldsymbol{\zeta}_j^T \boldsymbol{X}_i$, $\boldsymbol{\zeta}_j = \{\zeta_{j1}, \ldots, \zeta_{jq}\}^T$, and the $(j,k)$th element of $\boldsymbol{\Omega}_i$ is

$$\Omega_i(j,k) = \begin{cases} 1/\sigma_{\varepsilon_j}^2 & j = k \\ -\omega_{jk}(\boldsymbol{X}_i) & j > k \\ \Omega_i(k,j) & j < k \end{cases}.$$

In this model, the subject-dependent network connection between node $M_j$ and $M_k$ is modeled by a linear combination of $\boldsymbol{X}_i$ as $\omega_{jk}(\boldsymbol{X}_i) = \boldsymbol{\alpha}_{jk}^T \boldsymbol{X}_i$, where $\boldsymbol{\alpha}_{jk} = \{\alpha_{jk1}, \ldots, \alpha_{jkq}\}^T$. In the conditional Gaussian graphical model (1), we assume that the network structure is homogeneous in the whole population but network connection strengths vary across subjects and the heterogeneity in connection strengths depends on a vector of subject-specific covariates $\boldsymbol{X}_i$. When $\boldsymbol{\kappa}(\boldsymbol{X}_i)$ and $\Omega(\boldsymbol{X}_i)$ do not vary over $\boldsymbol{X}_i$, our model reduces to the regular Gaussian graphical models (Friedman *et al.*, 2008).

It is computationally intensive to directly maximize the joint likelihood in (1) that involves the unstructured precision matrix and it is difficult to calculate the gradients. Note that the number of parameters in the precision matrix is $p(p-1)/2 * q$. With $\boldsymbol{M}_i$ and $\boldsymbol{X}_i$ even at a moderate scale, the optimization will be over a high-dimensional parameter space: with $p = 20$ biomarkers and $q = 10$ covariates, the number of parameters is 2120. To compute the estimates efficiently, instead of maximizing the global likelihood function $P(\boldsymbol{M}_i|\boldsymbol{X}_i)$ in Equation (1), we optimize a pseudolikelihood formed by the products of all node-wise conditional likelihoods $P(M_{ij}|\boldsymbol{M}_{i,\backslash j}, \boldsymbol{X}_i)$, where $\boldsymbol{M}_{i,\backslash j}$ denotes a vector of $\boldsymbol{M}_i$ without the $j$th element. Replacing the global likelihood by pseudolikelihood results in simultaneously performing neighborhood estimation for all nodes while guaranteeing the symmetry of precision matrix, and provides consistent parameter estimates at the cost of statistical efficiency, but makes computation feasible in the presence of a large number of nodes and covariates. The conditional distribution of $M_{ij}$ given $\boldsymbol{M}_{i,\backslash j}$ and $\boldsymbol{X}_i$ is a normal distribution with a mean of $\mathbb{E}[M_{ij}|\boldsymbol{M}_{i,\backslash j}, \boldsymbol{X}_i] = \sigma_{\varepsilon_j}^2 (\boldsymbol{\zeta}_j^T \boldsymbol{X}_i + \sum_{k \neq j} \omega_{jk}(\boldsymbol{X}_i) M_{ik})$, and a variance of $\mathrm{Var}[M_{ij}|\boldsymbol{M}_{i,\backslash j}, \boldsymbol{X}_i] = \sigma_{\varepsilon_j}^2$. The pseudolikelihood function is then given by:

$$L_n(\boldsymbol{\zeta}, \boldsymbol{\alpha}, \sigma^2) = \prod_{i=1}^{n} \prod_{j=1}^{p} P(M_{ij}|\boldsymbol{M}_{i,\backslash j}, \boldsymbol{X}_i),$$

$$P(M_{ij}|\boldsymbol{M}_{i,\backslash j}, \boldsymbol{X}_i) = \sqrt{\frac{1}{2\pi\sigma_{\varepsilon_j}^2}}$$

$$\times \exp\left[-\frac{1}{2\sigma_{\varepsilon_j}^2}\left\{M_{ij} - \sigma_{\varepsilon_j}^2\left(\boldsymbol{\zeta}_j^T \boldsymbol{X}_i + \sum_{k \neq j} \omega_{jk}(\boldsymbol{X}_i) M_{ik}\right)\right\}^2\right].$$
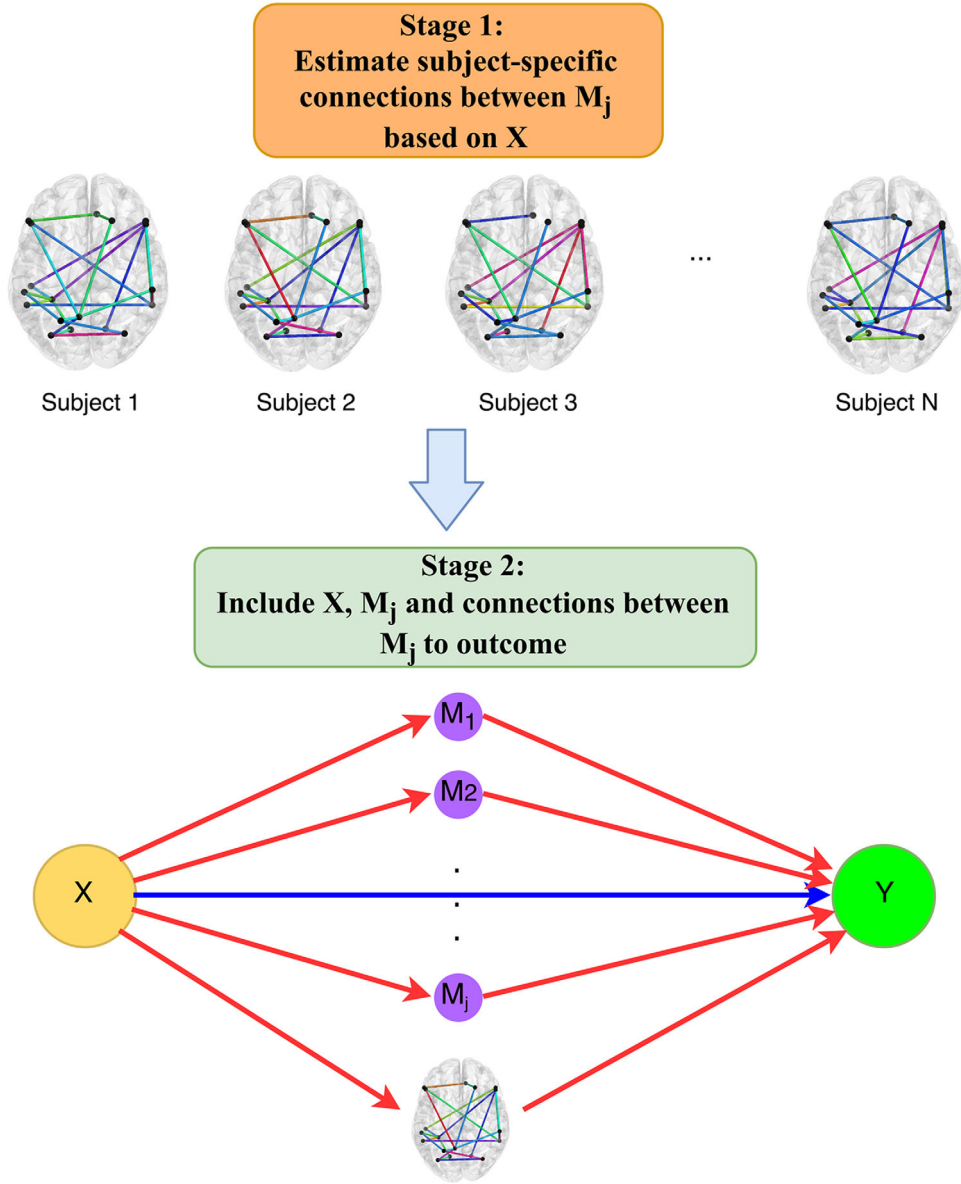
**FIGURE 1** Schematics of the proposed two-stage method

Denote the log-pseudolikelihood function by $l_n(\zeta, \alpha, \sigma^2) = \log L_n(\zeta, \alpha, \sigma^2)$. We impose regularization to stabilize estimation and perform variable selection. Our aim is to minimize the following penalized negative log-pseudo likelihood $-\frac{1}{n}l_n(\zeta, \alpha, \sigma^2) + p(\zeta, \alpha)$, where $p(\zeta, \alpha)$ is a properly chosen penalty function. The most desirable penalty for variable selection is to constrain the number of variables selected (eg, $l_0$-penalty), which is equivalent to the best subset selection. To this end, one may consider a penalty function $p(\zeta, \alpha) = \lambda_1(\sum_{j=1}^p \|\zeta_j\|_0 + \sum_{j>k} \|\alpha_{jk}\|_0)$, where $\|\zeta_j\|_0 = \sum_{s=1}^q I(\zeta_{js} \neq 0)$ and $\|\alpha_{jk}\|_0 = \sum_{s=1}^q I(\alpha_{jks} \neq 0)$. However, the optimization involving this penalty function is nondeterministic polynomial-time hard (NP-hard) due to the nonconvexity and discontinuity of the $l_0$-norm. A computation-

ally efficient two-step iterative procedure referred to as augmented penalized minimization-$L_0$ (APM-$L_0$) was proposed to approximate the $l_0$-penalization as closely as possible (Li *et al.*, 2018). APM-$L_0$ uses surrogate parameters to augment the $l_0$-norm of $(\zeta, \alpha)$ and guarantees that the surrogate parameters are close enough to $(\zeta, \alpha)$. It iteratively solves a penalized regression problem followed by a hard thresholding. In the first step of APM-$L_0$, we use a lasso penalty and employ an accelerated generalized gradient algorithm with backtracking for computation (Simon *et al.*, 2013).

In the second step of APM-$L_0$, hard-thresholding is performed to remove parameters with small magnitudes (Li *et al.*, 2018). We iteratively update between $\zeta_j, \alpha_{jk}$, and $\sigma_{\varepsilon_j}^2$ until convergence is reached. The tuning parameters are selected

using cross-validation to determine a final subject-dependent network model. The detailed update functions in the algorithm are presented in Web Appendix B.

## 2.2 | Second-stage clinical outcome model

In the first-stage model, our main interest is to estimate the biomarkers network effects (ie, $\omega_{jk}(X_i)$), as well as identifying important edges in the network. Thus, the first-stage method is an unsupervised approach to investigate basic brain organization through estimating subjects' network connections. In the second stage, we evaluate the clinical utility of the network connections identified in the first stage. We associate covariates, biomarkers, and the network connections with clinical impairment.

Under the conditional network model in (1), the heterogeneity in network connection across subjects is explained by the covariates. Mutual information, which quantifies the amount of information shared by two random variables, is used as a measure of network strength in many applications (eg, Song *et al.*, 2012). Under the jointly Gaussian assumption, the mutual information between two biomarkers given the remaining biomarkers is a transformation of the partial correlation, because the conditional joint distribution of the two biomarkers is a bivariate normal distribution. The transformation may alleviate the potential collinearity challenge between network connections when included in the outcome model. In addition, partial correlation yields an interpretation as conditional dependence and provides measures for determining modular membership in neuroimaging studies (Section 4.1, Web Appendix E2).

The second-stage model for disease outcome is:

$$E(Y_i|M_i, X_i) = \beta_0 + \beta^T X_i + \eta^T M_i + \sum_{s=1}^{p}\sum_{r=1}^{s-1}\gamma_{sr}W_{i,sr}, \quad (2)$$

where $W_{i,sr}$ is the mutual information defined as $W_{i,sr} = -\frac{1}{2}\log(1-\rho_{i,sr}^2)$, and $\rho_{i,sr} = -\frac{\Omega_i(s,r)}{\sqrt{\Omega_i(s,s)\Omega_i(r,r)}} = \sigma_{\varepsilon_s}\sigma_{\varepsilon_r}\omega_{sr}(X_i)$ is the partial correlation between $M_{i,s}$ and $M_{i,r}$ given $M_{i,\setminus\{s,r\}}$. The mutual information represents the strength of the connection between biomarkers given $X_i$. If $\rho_{i,sr} = 0$, then $W_{i,sr} = 0$, which implies that no connection is present between $M_{i,s}$ and $M_{i,r}$. If $\rho_{i,sr}$ is large, then $W_{i,sr}$ is large, which suggests a strong connection between $M_{i,s}$ and $M_{i,r}$.

In model (2), $\eta_s$, $\eta_r$, and $\gamma_{sr}$ are the effects of $X_i$ on $Y_i$ through network nodes and their connections, and $\beta$ are the effects of $X_i$ directly on $Y_i$ not through nodes and connections. Our interest is to evaluate the incremental effect of $\gamma_{sr}$ relative to $\eta_s$ and $\eta_r$. In our application, $\gamma_{sr}$ are the effects of causal gene on the clinical outcome through gray matter connection between a pair of ROIs ($M_{i,r}$, $M_{i,s}$).

We estimate parameters in the model (2) by minimizing a penalized least squares under the objective function $\frac{1}{2n}\sum_{i=1}^{n}\{Y_i - (\beta_0 + \beta^T X_i + \eta^T M_i + \sum_{s=1}^{p}\sum_{r=1}^{s-1}\gamma_{sr}\widehat{W}_{i,sr})\}^2 + q(\beta, \eta, \gamma)$, where $\widehat{W}_{i,sr} = -\frac{1}{2}\log(1-\widehat{\rho}_{i,sr}^2)$ are estimated based on the first-stage model, and the penalty function $q(\beta, \eta, \gamma) = \lambda_2(\|\beta\|_0 + \|\eta\|_0 + \|\gamma\|_0)$, with $\|\beta\|_0 = \sum_{s=1}^{q} I(\beta_s \neq 0)$, $\|\eta\|_0 = \sum_{j=1}^{p} I(\eta_j \neq 0)$, and $\|\gamma\|_0 = \sum_{s=1}^{p}\sum_{r=1}^{s-1} I(\gamma_{sr} \neq 0)$. We again use APM-$L_0$ with adaptive lasso penalty (the initial estimators are obtained from a ridge regression) on $\beta$, $\eta$, $\gamma$ to obtain parameter estimates. In practice, stability selection (Meinshausen and Bühlmann, 2010) can be further used to select informative network connections and important predictors of clinical outcomes. Stability selection combines subsampling and bootstrap with variable selection to provide improved performance (eg, reduced false discovery rate).

## 3 | SIMULATION STUDIES

We conducted extensive simulations to evaluate our method. We varied the number of covariates $q$, the number of biomarker nodes $p$, and the sample size $n$. Six settings were considered. The first four settings include $n = 500$ and $1000$ with $p = q = 5, 10, 20$; and $p > q$ with $(p, q) = (5, 3)$, $(p, q) = (10, 5)$, $(p, q) = (20, 10)$. In Setting 5, we considered $p = n = 100$ and $q = 5$. In Setting 6, we considered denser precision matrix and mean matrix with $p = 18, q = 17$, including a constant one in $X_i$ and $n = 500$ and $1000$.

In Settings 1 and 2, we simulated covariates $X_i$ independently from a standard normal distribution, where Setting 1 had homogeneous variance parameters $\sigma_{\varepsilon_j}^2 = 0.2$ and Setting 2 had heterogeneous variance parameters ($\sigma_{\varepsilon_j}^2 = 0.2$ or $0.15$). Four of $\alpha_{jk}$s are nonzeros with magnitudes ranging from $-1.5$ to $1.5$ and the remaining $\alpha_{jk}$s are all zeros. For example, $\alpha_{12} = \{-0.5, -1, -1.5, 0, 0\}^T$, $\alpha_{23} = \{-1, -0.5, -1.5, 0, 0\}^T$, $\alpha_{34} = \{1.5, -0.5, -1, 0, 0\}^T$, and $\alpha_{45} = \{-0.5, -1.5, 1, 0, 0\}^T$ when $p = q = 5$ in Setting 1. Setting 3 included an additional binary covariate and Setting 4 additionally examined scenarios with $p > q$. For Settings 1-5, we generated clinical outcomes in the second-stage model from $Y_i = X_{i1} + 2X_{i2} + M_{i1} + 3M_{i2} + W_{i,s_1r_1} + 2W_{i,s_5r_5} + \epsilon_i$, where $X_i$, $M_i$, and $W_i$ were standardized and $\epsilon_i \sim N(0, 1)$. Here, $W_{i,s_1r_1}$ is the mutual information between $M_{i1}$ and $M_{i2}$ for $p = 5, 10, 20$, and $W_{i,s_5r_5}$ is the mutual information between $M_{i2}$ and $M_{i3}$ for $p = 5$, and $M_{i1}$ and $M_{i6}$ for $p = 10, 20$. In Setting 6, 18 of $\alpha_{jk}$s are nonzeros and 14 of $\zeta_j$s are nonzeros. $\alpha_{jk}$s and $\zeta_j$s varied by different covariates. In addition, more network connections (10 connections) are associated with outcome. The details of the simulation settings are presented in Web Appendix C.

For each simulated data set, the length of the grid search vector of the tuning parameters $\lambda_1$ and $\lambda_2$ is 10 and 100,

**TABLE 1** Estimation and selection performance of simulations in Setting 1

| | | n = 500 | | | | | | n = 1000 | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | SSE[a] | TP[b] | FP[c] | TN[d] | FN[e] | MCC[f] | SSE[a] | TP[b] | FP[c] | TN[d] | FN[e] | MCC[f] |
| | | p = 5, q = 5 | | | | | | | | | | | |
| **First stage** | $\zeta$ | 0.409 | 8.97 | 1.44 | 14.56 | 0.03 | 0.901 | 0.167 | 9.00 | 0.94 | 15.06 | 0.00 | 0.931 |
| | $\alpha$ | 0.183 | 11.99 | 1.75 | 36.25 | 0.01 | 0.928 | 0.066 | 12.00 | 1.35 | 36.65 | 0.00 | 0.945 |
| **Second stage** | $\gamma$ | 0.050 | 2.00 | 0.39 | 7.61 | 0.00 | 0.917 | 0.020 | 2.00 | 0.26 | 7.74 | 0.00 | 0.940 |
| | $\beta$ | 0.006 | 2.00 | 0.22 | 2.78 | 0.00 | 0.916 | 0.003 | 2.00 | 0.32 | 2.68 | 0.00 | 0.907 |
| | $\eta$ | 0.021 | 2.00 | 0.36 | 2.64 | 0.00 | 0.891 | 0.007 | 2.00 | 0.25 | 2.75 | 0.00 | 0.920 |
| | | p = 10, q = 10 | | | | | | | | | | | |
| **First stage** | $\zeta$ | 1.106 | 8.89 | 3.10 | 87.90 | 0.11 | 0.872 | 0.455 | 9.00 | 2.03 | 88.97 | 0.00 | 0.925 |
| | $\alpha$ | 0.604 | 12.00 | 3.00 | 435.00 | 0.00 | 0.920 | 0.234 | 12.00 | 2.03 | 435.97 | 0.00 | 0.942 |
| **Second stage** | $\gamma$ | 0.067 | 2.00 | 0.43 | 42.57 | 0.00 | 0.922 | 0.023 | 2.00 | 0.36 | 42.64 | 0.00 | 0.936 |
| | $\beta$ | 0.012 | 2.00 | 1.06 | 6.94 | 0.00 | 0.812 | 0.004 | 2.00 | 0.74 | 7.26 | 0.00 | 0.854 |
| | $\eta$ | 0.029 | 2.00 | 1.16 | 6.84 | 0.00 | 0.789 | 0.011 | 2.00 | 0.95 | 7.05 | 0.00 | 0.826 |
| | | p = 20, q = 20 | | | | | | | | | | | |
| **First stage** | $\zeta$ | 2.369 | 8.43 | 2.84 | 388.16 | 0.57 | 0.871 | 0.969 | 8.98 | 3.53 | 387.47 | 0.02 | 0.903 |
| | $\alpha$ | 1.376 | 12.00 | 3.91 | 3784.09 | 0.00 | 0.903 | 0.622 | 12.00 | 3.08 | 3784.92 | 0.00 | 0.933 |
| **Second stage** | $\gamma$ | 0.068 | 2.00 | 0.33 | 187.67 | 0.00 | 0.943 | 0.030 | 2.00 | 0.24 | 187.76 | 0.00 | 0.958 |
| | $\beta$ | 0.012 | 2.00 | 1.85 | 16.15 | 0.00 | 0.765 | 0.007 | 2.00 | 2.04 | 15.96 | 0.00 | 0.756 |
| | $\eta$ | 0.031 | 2.00 | 1.98 | 16.02 | 0.00 | 0.751 | 0.013 | 2.00 | 2.05 | 15.95 | 0.00 | 0.739 |

[a]SSE: average sum of squared error across 100 simulations;
[b]TP: average number of true positive parameters across 100 simulations;
[c]FP: average number of false positive parameters across 100 simulations;
[d]TN: average number of true negative parameters across 100 simulations;
[e]FN: average number of false negative parameters across 100 simulations;
[f]MCC: Matthews correlation coefficient.

respectively. Ten-fold cross-validation was applied to choose the optimal tuning parameters. Simulations were repeated 100 times for each setting. To only retain the informative connectivities in predicting the disease outcome, we excluded near constant mutual information measures (variability<0.005) from the second-stage model.

Tables 1 and 2 and Web Tables S1, S2, S3, and S4 summarize the numeric simulation results in terms of sum of squared error (SSE), true positive (TP; number of nonnull variables correctly selected), false positive (FP; number of null variables incorrectly selected), true negative (TN; number of null variables correctly not selected), false negative (FN; number of nonnull variables incorrectly not selected) and Matthews correlation coefficient (MCC; Matthews, 1975), which is used as a measure that balances sensitivity and specificity even if the nonnull variable class and the null variable class are of very different size.

In all settings, our method yields small SSEs both in the first and second stages. In Setting 1 from $p = 5$ to $p = 20$ (Table 1), all of the TP parameters in $(\alpha_{jk})_{j>k}$ were selected, while small number of FP parameters were selected, and almost all of TN parameters were not selected in the first stage. MCCs of $(\alpha_{jk})_{j>k}$ were larger than 0.9. In the second stage,

our method identified all the TP variables in $W_i$, $X_i$, and $M_i$, along with a small number of FP variables and rejected most of TN variables. MCCs of $\gamma$ (connection effects) were larger than 0.9. Setting 2 is more difficult because $\sigma^2_{\varepsilon_j}$ varies by node. The performance in Setting 2 (Table 2) remained to be satisfactory with slightly more FPs and slightly lower MCC. When $p = q = 20$, the average number of FP parameters in $(\alpha_{jk})_{j>k}$ was 11.03 when $n = 500$ and it decreased to 3.91 when the sample size increased to 1000. Our method retained all TP variables in the second stage and the average number of FP variables in connections was still small. MCC of $\gamma$ was still beyond 0.85.

The performance in Settings 3-6 was similar. The results are presented in Web Appendix D. Web Figure S2 (Settings 1-4), Web Figure S3 (Setting 5), and Web Figure S4 (Setting 6) also visualize the number of times (at least once among 100 simulations) that an edge was identified in the network structure of $M_i$ in the first-stage analysis. Our method correctly identified all TP edges in all settings. In addition, we compared the estimated network in the first stage to the network estimated by the FGL in Setting 1 with $p = q = 20$. The FGL identified less true edges and selected many more null edges. The details are in Web Appendix D and Web Table S5.

**TABLE 2** Estimation and selection performance of simulations in setting 2

| | | n = 500 | | | | | | n = 1000 | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | SSE[a] | TP[b] | FP[c] | TN[d] | FN[e] | MCC[f] | SSE[a] | TP[b] | FP[c] | TN[d] | FN[e] | MCC[f] |
| | | $p = 5$, $\quad q = 5$ | | | | | | | | | | | |
| **First stage** | $\zeta$ | 0.425 | 8.98 | 1.33 | 14.67 | 0.02 | 0.906 | 0.176 | 9.00 | 1.09 | 14.91 | 0.00 | 0.922 |
| | $\alpha$ | 0.247 | 11.95 | 2.00 | 36.00 | 0.05 | 0.921 | 0.090 | 12.00 | 2.54 | 35.46 | 0.00 | 0.904 |
| **Second stage** | $\gamma$ | 0.020 | 2.00 | 0.51 | 7.49 | 0.00 | 0.889 | 0.006 | 2.00 | 0.52 | 7.48 | 0.00 | 0.883 |
| | $\beta$ | 0.007 | 2.00 | 0.39 | 2.61 | 0.00 | 0.892 | 0.003 | 2.00 | 0.42 | 2.58 | 0.00 | 0.884 |
| | $\eta$ | 0.041 | 2.00 | 0.50 | 2.50 | 0.00 | 0.873 | 0.016 | 2.00 | 0.42 | 2.58 | 0.00 | 0.874 |
| | | $p = 10$, $\quad q = 10$ | | | | | | | | | | | |
| **First stage** | $\zeta$ | 1.070 | 8.80 | 2.71 | 88.29 | 0.20 | 0.879 | 0.342 | 8.99 | 2.15 | 88.85 | 0.01 | 0.919 |
| | $\alpha$ | 0.691 | 11.99 | 2.50 | 435.5 | 0.01 | 0.928 | 0.243 | 12.00 | 1.89 | 436.11 | 0.00 | 0.945 |
| **Second stage** | $\gamma$ | 0.051 | 2.00 | 0.86 | 42.14 | 0.00 | 0.854 | 0.019 | 2.00 | 0.82 | 42.18 | 0.00 | 0.861 |
| | $\beta$ | 0.013 | 2.00 | 1.05 | 6.95 | 0.00 | 0.799 | 0.004 | 2.00 | 0.78 | 7.22 | 0.00 | 0.854 |
| | $\eta$ | 0.072 | 2.00 | 1.14 | 6.86 | 0.00 | 0.789 | 0.027 | 2.00 | 1.41 | 6.59 | 0.00 | 0.749 |
| | | $p = 20$, $\quad q = 20$ | | | | | | | | | | | |
| **First stage** | $\zeta$ | 2.605 | 8.59 | 3.52 | 387.48 | 0.41 | 0.867 | 1.000 | 9.00 | 2.65 | 388.35 | 0.00 | 0.917 |
| | $\alpha$ | 1.822 | 11.99 | 11.03 | 3776.97 | 0.01 | 0.808 | 0.836 | 12.00 | 3.91 | 3784.09 | 0.00 | 0.902 |
| **Second stage** | $\gamma$ | 0.036 | 2.00 | 0.73 | 187.27 | 0.00 | 0.876 | 0.014 | 2.00 | 0.81 | 187.19 | 0.00 | 0.867 |
| | $\beta$ | 0.014 | 2.00 | 2.02 | 15.98 | 0.00 | 0.758 | 0.006 | 2.00 | 2.17 | 15.83 | 0.00 | 0.726 |
| | $\eta$ | 0.041 | 2.00 | 2.43 | 15.57 | 0.00 | 0.706 | 0.017 | 2.00 | 2.54 | 15.46 | 0.00 | 0.698 |

[a]SSE: average sum of squared error across 100 simulations;

[b]TP: average number of true positive parameters across 100 simulations;

[c]FP: average number of false positive parameters across 100 simulations;

[d]TN: average number of true negative parameters across 100 simulations;

[e]FN: average number of false negative parameters across 100 simulations;

[f]MCC: Matthews correlation coefficient.

# 4 | APPLICATIONS TO THE GRAY MATTER NETWORK OF HD

HD is caused by an expansion of cytosine-adenine-guanine (CAG) triplet repeats in the *huntingtin* gene (MacDonald *et al.*, 1993). Existing studies (Johnson *et al.*, 2015) have shown that regional brain gray matter and white matter atrophy were associated with progression of HD (eg, the rate of change in patient's motor symptoms). In addition, McColgan *et al.* (2015) found that the degree (number of brain connections) of brain regions in white matter connectivity network was highly correlated with motor and cognitive deficits. These existing works considered brain subcortical volumes and cortical connectivity measures separately. Instead, here, we consider both types of measures on a subject.

We analyzed data collected from the recently completed PREDICT-HD (Paulsen *et al.*, 2014), a long-term natural history study of premanifest HD gene-positive subjects. Our analysis consisted of 499 subjects who carried an expanded CAG repeats in the premanifest stage without an HD diagnosis at the baseline visit. The median follow-up length was 6.3 years. In our analyses, $\boldsymbol{M}_i$ includes baseline brain atrophy measures obtained from structural MRI. The regional subcortical volumetric measures and cortical thickness measures were preprocessed using Freesurfer 5.3 and the details were described in Paulsen *et al.* (2014). All volumetric measures were adjusted for total intracranial volume. Six subcortical ROI gray matter volumes and 18 cortical thickness ROIs were included in the analyses (see Web Table S6). These ROIs were selected from a marginal screening (linear regression with second-stage clinical measure as outcome) based on false discovery-rate-corrected *p*-values.

The covariates $\boldsymbol{X}_i$ include 15 baseline clinical variables (details in Web Appendix E1), total gray matter volume, and total white matter volume. The second-stage model outcome was the rate of change in total motor score (TMS) estimated by a linear mixed-effects model with subject-specific random intercepts and random slopes, treating time since baseline as the time scale and adjusting for the baseline TMS.

In the first-stage analyses of PREDICT-HD, we estimated the network structure and connection strength of cortical thickness and subcortical volumes separately under model (1). In the second stage, we included covariates, subcortical volumes, cortical thickness, and the subcortical and cortical
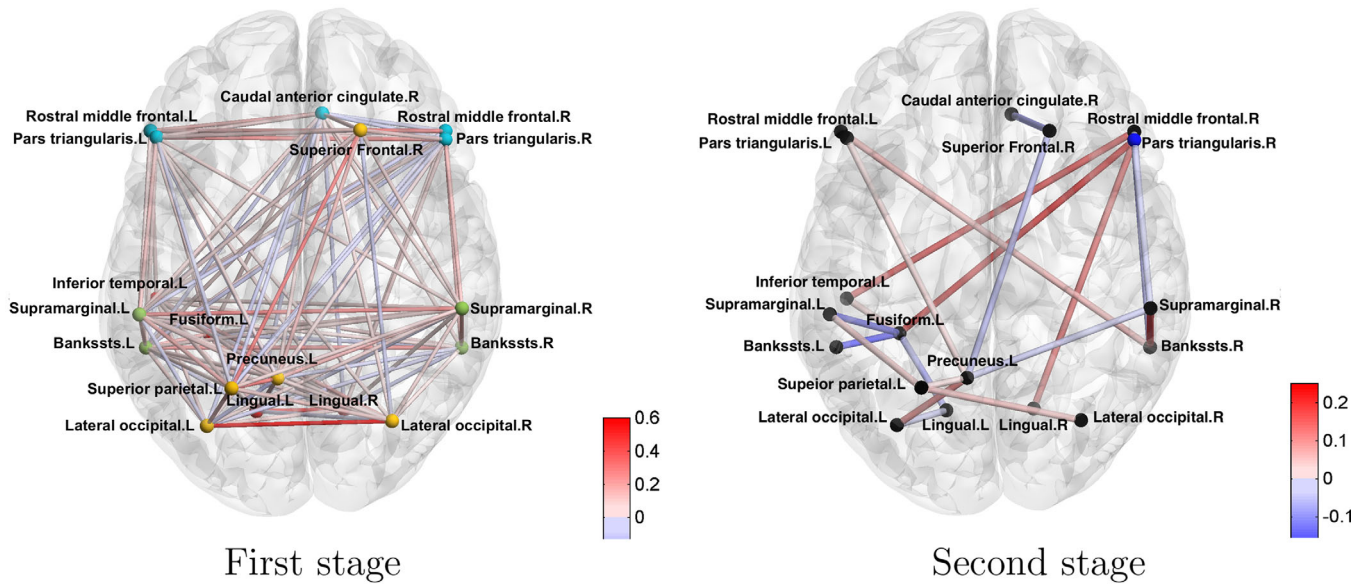
**FIGURE 2** Estimated cortical network (first stage) and effects of cortical connections on the rate of change of the total motor score (second stage). In the left panel, edge color represents the estimated partial correlation for an average subject (covariates fixed at the sample averages). Color of nodes represents modular membership. In the right panel, edge color represents the estimated effect size of the cortical connection on the rate of change of the total motor score. Blue nodes: cortical thickness ROIs that are also selected in the second-stage model. Black nodes: cortical thickness ROIs themselves are not selected in the second stage model, but their connections are selected

network connections obtained from the first stage to identify the important predictors of the rate of change in TMS under model (2). Stability selection (Meinshausen and Bühlmann, 2010) with 100 bootstraps was used to select informative network connections in the first stage and important predictors in the second stage. For each bootstrap sample, 10-fold cross-validation was applied to choose the optimal tuning parameter $\lambda_1$ in the first stage. The selection rule of the parameters in $(\zeta_j)$ and $(\alpha_{jk})_{j>k}$ was based on the relative frequency among the 100 times bootstrap. A variable would be selected if its relative frequency was larger than or equal to a threshold. To reduce the risk of missing important connections in the first-stage analysis, the threshold was set at 0.5 (Meinshausen and Bühlmann, 2010). In the second stage, we constructed a set of 100 initial regularization parameters, denoted by $\Lambda_2$, and performed variable selection on each $\lambda_2 \in \Lambda_2$. Hard-thresholding was implemented and the total number of selected variables was set to be 100. The final model was determined by the stability selection with a threshold of 0.65. We refitted the models on the full data including only the identified predictors and network connections to provide final parameter estimates.

## 4.1 | First-stage analysis results

Among 270 potential parameters in the subcortical network, 150 were selected to be informative (nonzero). For the cortical network, 908 out of 2754 parameters were nonnull. The average cortical network identified in the first stage is visualized in Figure 2 (left panel). In this figure, the strength of the edge between two ROIs represents the estimated partial correlation between them for an "average" subject with covariates fixed at the sample averages. Two strongest connections are the interhemispheric links between the left and right lateral occipital regions (contralateral homologous regions) and between the left superior parietal and right superior frontal regions. The degrees (defined as the number of links between a target ROI and other ROIs) of the cortical nodes for an average subject are visualized in Figure 3. The largest degree (degree = 17) was seen in the right caudal anterior cingulate and left bankssts regions, whereas the smallest degree (degree = 11) appeared in the left pars triangularis and left fusiform regions.

To understand the network organization of the cortical ROIs, they were classified into four distinct modules using the community Louvain algorithm based on the modularity optimization (Blondel *et al.*, 2008). Our network modularity is consistent with existing literature and the details are described in Web Appendix E2. From the first-stage analysis, our conditional Gaussian graphical model reveals that the cortical thickness network is dense but organizes in a modular fashion in premanifest HD, so that brain regional cortical thinning acts in a dependent manner. In Web Appendix E3, we present the identified top covariates that are associated with cortical connections in the first stage.
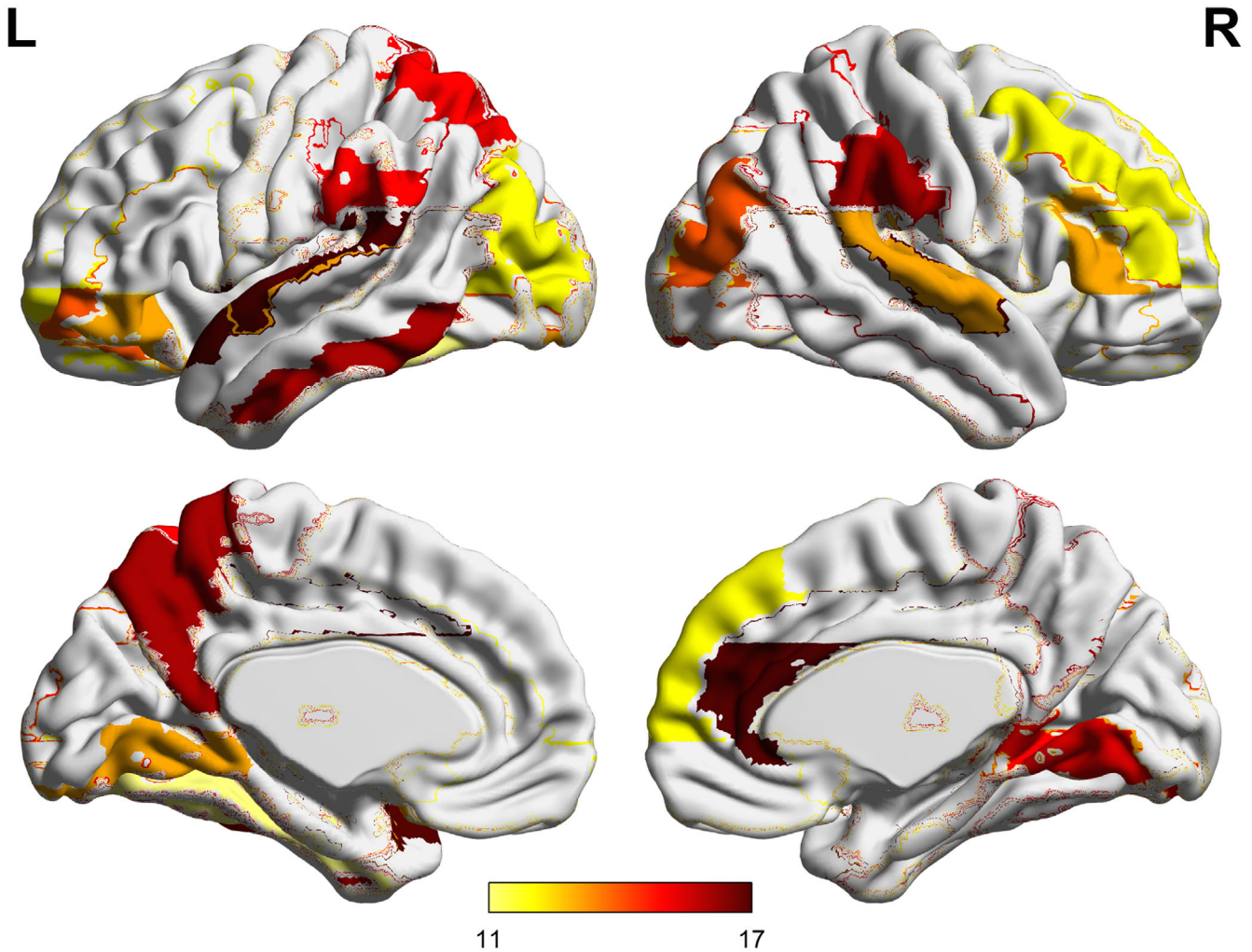
**FIGURE 3** Degree of cortical nodes selected in the first stage for an average subject. Top row: lateral view of left and right hemispheres. Bottom row: medial view of left and right hemispheres. Color represents degree

## 4.2 | Second-stage analysis results

Among all covariates $X_i$, CAP score and white matter total volumes are identified as informative, while none of the cognitive or functioning measures is identified given other variables in the final model. Among subcortical measures, the volumetric ROIs are important for predicting TMS rate of change, but not their network connections (only one connection was identified). In contrast, cortical thickness network connections, instead of cortical thickness ROI measures, are identified as important. Caudate, globus pallidus, and thalamus are selected among six subcortical volumes, whereas right pars triangularis is the only one selected among 18 cortical thickness ROIs. The connection between thalamus and globus pallidus regions is the only subcortical connection identified, whereas 22 out of 153 cortical connections (Figure 2, right panel) are identified to be predictive of motor symptom change.

The brain gray matter cortical thickness network identified in the first stage is dense (126 out of 153 connections were identified for a population average subject). However, most of connections do not predict TMS rate of change, which results in a sparse model in the second stage. The effect size of each cortical connection is summarized in Web Table S8.

Five cortical connections mostly predictive of the motor symptom are the interhemispheric connections (Web Table S8). Several anatomically long-distance connections are also identified, for example, the connection between right lingual and right pars triangularis. In addition, most of the connections with positive effects are the interhemispheric connections or anatomically long-distance connections. Web Figure S5 shows the identified paths from a baseline covariate to the clinical motor symptom through brain gray matter connections and our results were consistent with clinical literature. The details are in Web Appendix E5.

To assess the incremental predictive power of the selected ROIs and connections, we randomly selected two-thirds of subjects as a training set and the remaining one-third subjects as testing data and compared the $R$-squared value of the best performing model reported in the literature without imaging measures (with covariates TMS, symbol digit modality test raw score (SDMT), and CAP (a product of CAG and age measuring disease burden); Long *et al.*, 2017) with our model including imaging biomarkers and their connections (ie, TMS, SDMT, CAP, regional subcortical volumetric and cortical thickness $M_i$, subcortical and cortical connections $W_i$). We repeated this process 100 times. A ridge-penalty was imposed on cortical connections to minimize overfitting, except for the connections between right superior frontal and left precuneus regions and between left superior parietal and left precuneus regions, which were identified in the previous literature (Chen *et al.*, 2011). The average variance explained by the imaging biomarker model with cortical connections was 30.1%, as compared to 25.7% of the standard nonimaging model. In Web Appendix E6, we also showed that the imaging biomarker model improves the net reclassification rate for predicting conversion to HD diagnosis.

## 4.3 | Validation on TRACK-HD

We sought to validate findings from PREDICT-HD study on an independent natural history study of HD, TRACK-HD (Tabrizi *et al.*, 2009), which collected comprehensive gray matter and white matter structural neuroimaging measures and clinical assessments of premanifest HD patients. The cohort in the replication analyses includes 96 premanifest HD subjects with CAG repeat expansions. The cohort was followed up at four time points (year 2008, 2009, 2010, and 2011) and cortical thickness ROIs were obtained from structural MRI and preprocessed by Freesurfer in a similar fashion as PREDICT-HD (McColgan *et al.*, 2015, 2017).

We evaluated the predictive performance of 22 cortical connections identified to be informative in PREDICT on an independent study, TRACK-HD. First, we recalibrated strength of the 22 cortical connections identified in PREDICT on TRACK. Not all the measures used in the first-stage model of the PREDICT analyses were available for TRACK. Thus, the subject-specific cortical connections are estimated using baseline covariates (details in Web Appendix F). Similar to the PREDICT analyses, we compared the leave-one-out $R$-squared value and mean squared error (MSE) of two linear models predicting TMS at the last visit: a standard nonimaging model including baseline TMS and CAP score at the last visit as covariates, and a brain gray matter imaging biomarker model additionally including the 22 estimated cortical connections. A ridge penalty was imposed on cortical

connections to minimize overfitting, which is the same as PREDICT analyses. Leave-one-out $R$-squared value of the imaging model was 35.3% compared to the nonimaging model of 26.7%. Thus, on the independent TRACK-HD study, imaging biomarkers explained 0.32-fold additional variance of the TMS. In addition, the decrease in the MSE of the imaging model was 13% (from 0.725 point/year to 0.641 point/year). The results show that the identified cortical connections are useful in predicting follow-up TMS in an independent HD study.

We also compared our method with FGL. The cohort was first classified into a high-risk and low-risk classes based on median split of the CAP score. The 22 cortical connections were recalibrated through FGL with tuning parameters selected by Akaike information criterion (AIC). The leave-one-out $R$-squared value of the imaging biomarker model was only 28.1% and the MSE was 0.711 point/year, which suggests that using our method to estimate subject-specific connections explains more variance in TMS.

## 4.4 | Biological implications and insights

We compared the obtained gray matter networks with the white matter connectivity networks obtained on the TRACK-HD (McColgan *et al.*, 2017). Highly similar patterns are observed between our gray matter connection study and the white matter connectivity study, which suggest a shared underlying mechanism of HD. These results address an important biological question of whether white matter loss is a direct result of neuronal loss or loss of myelin or dysmyelination. The shared patterns support the former hypothesis. The details are in Web Appendix G.

## 5 | DISCUSSION

In this work, we propose a two-stage method to estimate a subject-dependent network from conditional Gaussian graphical model and evaluate the incremental effect of the network measures on a clinical outcome. Our method can simultaneously handle biomarkers and the between-subject heterogeneity of the biomarker network measures using covariates. Our analyses results suggest that brain cortical gray matter network connections are predictive of HD motor impairment in addition to regional atrophy and other HD risk factors.

Several extensions can be considered. Similar to neighborhood-based methods for graph selection (Meinshausen and Bühlmann, 2006; Peng *et al.*, 2009), our method is computationally attractive but cannot guarantee precision matrix to be positive-definite. We do not include such a constraint because our goal is to extract useful features from the network to predict clinical outcomes instead

of estimating the joint distribution of nodes. To ensure the positive definiteness of the estimated precision matrix, one can replace current neighborhood-based method with global likelihood-based approach along the line of graphical lasso that automatically guarantees positive-definite (Friedman *et al.*, 2008; Tibshirani *et al.*, 2015) or impose a positive-definite constraint to our algorithm, but at the cost of computational complexity. In the subject-dependent network model (1), we assume that network structure is homogeneous in the whole population because in our application, subjects were recruited at a similar clinical stage (premanifest stage). For other applications, to obtain subject-specific network structure, we can further extend our method to select network edges based on subject-specific connection strength.

## ORCID
*Shanghong Xie* https://orcid.org/0000-0001-5132-3887
*Donglin Zeng* https://orcid.org/0000-0003-0843-9280
*Yuanjia Wang* https://orcid.org/0000-0002-1510-3315

## REFERENCES

Alexander-Bloch, A., Giedd, J.N. and Bullmore, E. (2013) Imaging structural co-variance between human brain regions. *Nature Reviews Neuroscience*, 14, 322–336.

Alexander-Bloch, A., Raznahan, A., Bullmore, E. and Giedd, J. (2013) The convergence of maturational change and structural covariance in human cortical networks. *Journal of Neuroscience*, 33, 2889–2899.

Blondel, V.D., Guillaume, J.-L., Lambiotte, R. and Lefebvre, E. (2008) Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2008, P10008.

Cai, T.T., Li, H., Liu, W. and Xie, J. (2012) Covariate-adjusted precision matrix estimation with an application in genetical genomics. *Biometrika*, 100, 139–156.

Chen, Z.J., He, Y., Rosa-Neto, P., Gong, G. and Evans, A.C. (2011) Age-related alterations in the modular organization of structural cortical network by using cortical thickness from mri. *Neuroimage*, 56, 235–245.

Chen, M., Ren, Z., Zhao, H. and Zhou, H. (2016) Asymptotically normal and efficient estimation of covariate-adjusted Gaussian graphical model. *Journal of the American Statistical Association*, 111, 394–406.

Cheng, J., Levina, E., Wang, P. and Zhu, J. (2014) A sparse Ising model with covariates. *Biometrics*, 70, 943–953.

Danaher, P., Wang, P. and Witten, D.M. (2014) The joint graphical Lasso for inverse covariance estimation across multiple classes. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76, 373–397.

Friedman, J., Hastie, T. and Tibshirani, R. (2008) Sparse inverse covariance estimation with the graphical Lasso. *Biostatistics*, 9, 432–441.

He, Y., Chen, Z. and Evans, A. (2008) Structural insights into aberrant topological patterns of large-scale cortical networks in Alzheimer's disease. *Journal of Neuroscience*, 28, 4756–4766.

Johnson, E.B., Rees, E.M., Labuschagne, I., Durr, A., Leavitt, B.R., Roos, R.A., Reilmann, R., Johnson, H., Hobbs, N.Z., Langbehn, D.R., Stout, J.C., Tabrizi, S.J., Scahill, R.I. and TRACK-HD Investigators. (2015) The impact of occipital lobe cortical thickness on cognitive task performance: an investigation in Huntington's disease. *Neuropsychologia*, 79, 138–146.

Li, X., Xie, S., Zeng, D. and Wang, Y. (2018) Efficient $l_0$-norm feature selection based on augmented and penalized minimization. *Statistics in Medicine*, 37, 473–486.

Long, J.D., Langbehn, D.R., Tabrizi, S.J., Landwehrmeyer, B.G., Paulsen, J.S., Warner, J. and Sampaio, C. (2017) Validation of a prognostic index for Huntington's disease. *Movement Disorders*, 32, 256–263.

MacDonald, M.E., Ambrose, C.M., Duyao, M.P., Myers, R.H., Lin, C., Srinidhi, L., Barnes, G., Taylor, S.A., James, M., Groot, N., MacFarlance, H., Jenkins, B., Anderson, M.A., Wexler, N.S. and Gusella, J.F. (1993) A novel gene containing a trinucleotide repeat that is expanded and unstable on Huntington's disease chromosomes. *Cell*, 72, 971–983.

Matthews, B.W. (1975) Comparison of the predicted and observed secondary structure of t4 phage lysozyme. *Biochimica et Biophysica Acta (BBA)-Protein Structure*, 405, 442–451.

McColgan, P., Seunarine, K.K., Gregory, S., Razi, A., Papoutsi, M., Long, J.D., Mills, J.A., Johnson, E., Durr, A., Roos, R.A., Leavitt, B.R., Stout, J.C., Scahill, R.I., Clark, C.A., Rees, G., Tabrizi, S.J. and Track-On HD Investigators. (2017) Topological length of white matter connections predicts their rate of atrophy in premanifest Huntington's disease. *JCI Insight*, 2, e92641.

McColgan, P., Seunarine, K.K., Razi, A., Cole, J.H., Gregory, S., Durr, A., Roos, R.A., Stout, J.C., Landwehrmeyer, B., Scahill, R.I., Clark, C.A., Rees, G., Tabrizi, S.J. and Track-HD Investigators. (2015) Selective vulnerability of rich club brain regions is an organizational principle of structural connectivity loss in Huntington's disease. *Brain*, 138, 3327–3344.

Meinshausen, N. and Bühlmann, P. (2006) High-dimensional graphs and variable selection with the Lasso. *The Annals of Statistics*, 34, 1436–1462.

Meinshausen, N. and Bühlmann, P. (2010) Stability selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72, 417–473.

Paulsen, J.S., Long, J.D., Johnson, H.J., Aylward, E.H., Ross, C.A., Williams, J.K., Nance, M.A., Erwin, C.J., Westervelt, H.J., Harrington, D.L., Bockholt, H.J., Zhang, Y., McCusker, E.A., Chiu, E.M., Panegyres, P.K. and PREDICT-HD Investigators and Coordinators of the Huntington Study Group. (2014) Clinical and biomarker changes in premanifest Huntington disease show trial feasibility: a decade of the PREDICT-HD study. *Frontiers in Aging Neuroscience*, 6, 78.

Peng, J., Wang, P., Zhou, N. and Zhu, J. (2009) Partial correlation estimation by joint sparse regression models. *Journal of the American Statistical Association*, 104, 735–746.

Simon, N., Friedman, J., Hastie, T. and Tibshirani, R. (2013) A sparse-group Lasso. *Journal of Computational and Graphical Statistics*, 22, 231–245.

Song, L., Langfelder, P. and Horvath, S. (2012) Comparison of co-expression measures: mutual information, correlation, and model based indices. *BMC Bioinformatics*, 13, 328.

Tabrizi, S.J., Langbehn, D.R., Leavitt, B.R., Roos, R.A., Durr, A., Craufurd, D., Kennard, C., Hicks, S.L., Fox, N.C., Scahill, R.I., Borowsky, B., Tobin, A.J., Rosas, H.D., Johnson, H., Reilmann, R., Landwehrmeyer, B., Stout, J.C. and TRACK-HD Investigators. (2009) Biological and clinical manifestations of Huntington's disease in the longitudinal TRACK-HD study: cross-sectional analysis of baseline data. *The Lancet Neurology*, 8, 791–801.

Yin, J. and Li, H. (2011) A sparse conditional Gaussian graphical model for analysis of genetical genomics data. *The Annals of Applied Statistics*, 5, 2630.

## SUPPORTING INFORMATION

Web appendices, tables, and figures referenced in Sections 1-4 and code are available with this paper at the Biometrics website on Wiley Online Library.

**How to cite this article:** Xie S, Li X, McColgan P, Scahill RI, Zeng D, Wang Y. Identifying disease-associated biomarker network features through conditional graphical model. *Biometrics*. 2019;1–12. https://doi.org/10.1111/biom.13201