1  # Ensemble-based network aggregation improves the accuracy of
2  # gene network reconstruction

3  Jeffrey D. Allen[1,2], Yang Xie[1,2] and Guanghua Xiao[1,*]

4  [1] Quantitative Biomedical Research Center
5  [2] Simmons Cancer Center, UT Southwestern Medical Center,

6  * To whom correspondence should be addressed. Tel: 214-648-4553; Fax: 214-648-7673; Email:
7  Guanghua.Xiao@UTSouthwestern.edu

8  Present Address:  Guanghua Xiao, 5323 Harry Hines Blvd, UT Southwestern Medical Center, Dallas,
9  TX, 75390, USA

10

11  **ABSTRACT**

12  Reverse engineering approaches to construct context-specific gene regulatory networks (GRNs)
13  based on genome-wide mRNA expression data have led to significant biological findings. However,
14  the reliability and reproducibility of the reconstructed GRNs needs to be improved. Here, we propose
15  an ensemble-based network aggregation approach to improve the accuracy of the network topology
16  constructed from mRNA expression data. To evaluate the performance of different approaches, we
17  created dozens of simulated networks and also tested our methods on three Escherichia coli datasets.
18  We demonstrate three novel applications from this development. First, bootstrapping can be done on
19  the available samples, turning any network reconstruction approach into an ensemble method.
20  Second, this aggregation approach can be used to combine GRNs from different network inference
21  methods, creating a novel network reconstruction approach that consistently outperforms any
22  constituent method. Third, the approach can be used to effectively integrate GRNs constructed from
23  different studies – producing more accurate networks. We are releasing an implementation of these
24  techniques as an R package "ENA" which is able to run network inference in parallel across multiple
25  servers. We made all of the code and data used in our simulations and analysis available online at
26  https://github.com/QBRC/ENA-Research to ensure the reproducibility of our results.

27  **INTRODUCTION**

28      Accurate reconstruction of Gene Regulatory Networks (GRNs) from gene expression

29  microarrays has been shown to be valuable in a myriad of areas surrounding biomedical research[1–

30  5]. Researchers have previously used approaches including Bayesian Network-Based approaches [6],

31  Correlation-Based approaches [7], and Partial-Correlation-Based approaches [8,9]. These methods

32  have been shown to have various strengths and weaknesses under different biological/simulation

33  settings, with no one method excelling in all conditions. [10]. Additionally, leveraging gene expression

34  data from multiple datasets to construct gene networks is often difficult, due to discrepancies in

35  microarray platform selection, as well as in normalization and data processing techniques. In this

36  study, we propose an Ensemble-based Network Aggregation (ENA) approach to integrate gene

37  networks derived from different methods and different datasets to improve the accuracy of network

38  inference.

39      We used a non-parametric, inverse-rank-product, algorithm in the ENA approach to combine

40  networks reconstructed on the same set of genes. The rank- product method was introduced by

41  Breitling et al [11,12] as an effective method for detecting differentially expressed genes in microarray

42  studies. Because the rank product method is powerful and computationally efficient, it was extended

43  to be used in other fields, such as RNAi screening [13] and proteomics [14]. This method can be

44  directly related to the linear rank statistics [15]. In this study, we show three ways to leverage this

45  approach to generate the ensemble-based networks: 1.) Samples in a dataset can be "bootstrapped"

46  to reconstruct multiple networks out of a single original dataset using a single reconstruction method,

47  which can then be aggregated into a more accurate and reproducible network; 2.) Networks produced

48  by various reconstruction methods can be aggregated into a single network that is more accurate than

49  the network provided by any individual method; 3.) Networks reconstructed from different studies

50  which contain the same genes can be combined into a single, more accurate network, despite

51  differences in platforms or normalization techniques. Because this approach has little overhead, it can

52  efficiently be applied to dozens or hundreds of networks reconstructed on the same set of genes. We

53  find that this approach has the ability to improve the accuracy of GRN reconstruction in all three

54  applications based on simulated gene expression data, as well as *Escherichia coli* (E. coli)

55  datasets[16–19].

56  **MATERIAL AND METHODS**

57  **Overview of the Inverse-Rank-Product Network Aggregation Approach**

58  Reconstructed gene networks are often returned as a weighted undirected graph $G = (N, \Omega)$, where

59  $G$ is a reconstructed graph, $N = \{1, ..., n\}$ is the set of vertices (genes) in the graph, and

60  $\Omega = [\omega_{ij}]_{i, j \in N}$ is referred to as the adjacency matrix, in which $\omega_{ij}$ represents the intensity of the

61  interaction between genes i and j. A larger (absolute) value of $\omega_{ij}$ indicates a stronger interaction or

62  higher confidence in the edge between genes i and j, while $\omega_{ij} = 0$ indicates no interaction, or

63  conditional independence between genes i and j. Some techniques, such as Sparse PArtial

64  Correlation Estimations (SPACE) [9], return a sparse matrix in which many of the possible interactions

65  are 0; other techniques return complete graphs in which all edges are present with non-zero

66  weightings. Additionally, the distribution of $\omega_i$ can vary drastically among reconstruction techniques.

67  For this reason, the aggregation of networks reconstructed using different techniques or different

68  datasets is challenging. In this study, we used a rank-product method to combine networks to

69  overcome the different distributions observed in this problem.

70  Specifically, suppose $\boldsymbol{G} = \{G^k\}$ is a set of networks constructed on the same set of genes N,

71  where $k = \{1, ..., K\}$ is the index of a particular network. For each single network $G^k = (N, \Omega^k)$, we

72  calculate $r_{ij}^k$, the rank of $\omega_{ij}^k$ for $\{i, j \in N \text{ and } i < j\}$. Since the adjacency matrix $\Omega$ of an

73  undirected graph is a symmetric matrix, we only need to calculate the rank of the $N*(N-1)/2$

74  elements in $\omega_{ij}$ constituting the lower triangle (i < j) of the adjacency matrix. In this study, we give the

75  lower rank to the high strength/confidence interaction. For example, the interaction with the highest

76  strength/confidence will have rank 1. This operation is performed on each individual graph $G^k$

77  independently. After the rank of $r_{ij}^k$ has been computed for each network $G^k$, we calculate the rank

78  of a particular edge between genes i and j in the aggregated network by taking the product of the

79  ranks of the same edge across all networks in $\boldsymbol{G}$, as follows: $\tilde{r}_{ij} = \prod_{k=1}^{K} r_{ij}^k$ .This function is iterated

80  over all possible edges to construct the aggregated network $\tilde{G} = (N, \tilde{r}_{ij})$, in which the strength of the

81  edges in the new network are based on the aforementioned rank-product calculation.

82        This algorithm can be efficiently applied to large networks with many reconstructed networks

83  in $\boldsymbol{G}$. The complexity of the algorithm is $O(K \cdot |N| \log(|N|))$, as $\dfrac{|N|^2 - |N|}{2} = O(N^2)$ elements must

84  be sorted for each network in $G^k$.

**Three Applications for Ensemble-based Network Aggregation**

86        The initial application was to leverage the rank-product method to "bootstrap" samples. Each

87  time, we construct the gene network using a randomly selected subset of the available samples. By

88  repeating this process B times, we create a set $\boldsymbol{G}$ consisting of B graphs, each reconstructed using

89  only randomly selected bootstrap samples in the dataset. For example, here is the procedure to

90  generate the bootstrapping network from the microarray dataset MD:

$$MD \xrightarrow{Bootstrap} \begin{cases} MD^1 & \rightarrow & G^1 = \{N, \Omega^1\} & \rightarrow & r_{ij}^1 \text{ (for } 1 \le i < j \le n) \\ \vdots & & \vdots & & \vdots & \rightarrow \text{RankProduct} \rightarrow \tilde{G} \\ MD^B & \rightarrow & G^B = \{N, \Omega^B\} & \rightarrow & r_{ij}^B \text{ (for } 1 \le i < j \le n) \end{cases}$$

91
92

93        Of course, this bootstrapping procedure inflates the computational complexity of GRN

94  reconstruction by orders of magnitude, as GRNs must be reconstructed B times, rather than just once.

95  Because each graph in $\boldsymbol{G}$ can be reconstructed independently, it is possible to take advantage of the

96  "parallelizability" of these simulations by utilizing multiple cores or computers as we discuss below.

97  Note also that the complexity of GRN reconstruction does scale on the order of samples included, so

98  each permuted GRN can be constructed slightly more quickly than a single global GRN; for the

99 reconstruction techniques employed in this study, however, the performance did not vary greatly

100 based on the number of samples included.

101 The second application of the rank-product network merging method was to reconstruct an

102 aggregated GRN based on the output of multiple different reconstruction techniques. We have

103 observed that reconstruction techniques perform differently based on different simulation settings [20],

104 with no one method outperforming the others on all metrics. Thus, we were interested to see whether

105 or not merging these GRNs would offer an improvement in performance. In this application, the set of

106 graphs $G$ consist of one graph per network reconstruction technique employed. In our analysis, we

107 leveraged GeneNet[8], Weighted Correlation Network Analysis (WGCNA) [7], and Sparse PArtial

108 Correlation Estimation (SPACE) [9], creating a set of 3 graphs which can then be aggregated.

109 GeneNet and SPACE are partial-correlation-based inference algorithms. GeneNet uses the Moore-

110 Penrose pseudoinverse [21] and bootstrapping to estimate the concentration matrix. The SPACE

111 algorithm creates a regression problem when trying to estimate the concentration matrix and then

112 optimizes the results with a symmetric constraint and an L1 penalization. WGCNA is a correlation-

113 based approach which can identify sub-networks using hierarchical clustering. Conceptually, the

114 aggregated graph would place higher confidence on those edges which were consistently ranked

115 highly across the three methods, and would place lower confidence on those edges which were only

116 ranked highly in one graph. This is the procedure to derive the ensemble network based on M

117 different methods on the same dataset MD:

$$
MD \begin{cases}
\xrightarrow{\text{method 1}} & G^1 = \{N, \Omega^1\} \quad \rightarrow \quad r_{ij}^1 \text{ (for } 1 \le i < j \le n) \\
\vdots & \vdots \qquad\qquad\qquad \vdots \qquad\qquad \rightarrow \text{RankProduct} \rightarrow \tilde{G} \\
\xrightarrow{\text{method M}} & G^M = \{N, \Omega^M\} \quad \rightarrow \quad r_{ij}^M \text{ (for } 1 \le i < j \le n)
\end{cases}
$$

118

119 The final application evaluated in this study was in the merging of networks constructed from

120 different datasets. Historically, gene expression datasets have been collected from various sites on

121 different microarray platforms with different procedures for tissue collection; this creates

122 incompatibilities and difficulties when trying to perform analysis on data from different datasets

123 simultaneously. Because the rank-product method makes no assumptions on the distribution of the

124 data at any point, we employ it to combine GRNs produced from different datasets, yielding a single,

125 aggregated GRN which aims to capture the consistencies in network topology from the GRNs

126 produced on different datasets. Here is the procedure to derive the aggregated network from datasets

127 $MD^1$, $MD^2$…. $MD^D$:

$$
\begin{aligned}
MD^1 &\rightarrow \quad G^1 = \{N, \Omega^1\} \quad \rightarrow \quad r_{ij}^1 \text{ (for } 1 \le i < j \le n) \\
\vdots &\qquad \vdots \qquad\qquad\qquad\qquad \vdots \qquad\qquad \rightarrow \text{RankProduct} \rightarrow \tilde{G} \\
MD^D &\rightarrow \quad G^D = \{N, \Omega^D\} \quad \rightarrow \quad r_{ij}^D \text{ (for } 1 \le i < j \le n)
\end{aligned}
$$

128

129 **Software**

130     The code used to bootstrap samples and aggregate the resultant networks was written in the

131     R programming language [22]. We created an R Package entitled "ENA" and have made it available

132     on CRAN (http://cran.r-project.org/web/packages/ENA/index.html); the compiled binaries, as well as

133     all original source code are available for download there.

134     Because of the parallelization opportunities in this algorithm, we ensured that our software

135     would be able to distribute the bootstrapping process across multiple cores and multiple nodes using

136     MPI [23]. Thus, if 150 CPU cores were available simultaneously, a bootstrapping of 150 samples

137     could run in approximately the same amount of wall-clock time as a single reconstruction using all

138     samples could. The ENA package includes robust documentation and (optionally) leverages the RMPI

139     package to allow parallel execution of the bootstrapping simulations where such a computational

140     infrastructure is available.

141     Additionally, we leveraged the Git revision control system via GitHub (http://github.com) to

142     control not only the R code developed for the ENA package, but also all code, reports, and data used

143     in the aforementioned simulations and reconstruction techniques; all of this code is freely available at

144     https://github.com/QBRC/ENA-Research. All the data analysis code that has been used to generate

145     the results in this study was compiled into a single report and can be reproduced easily by using the

146     knitr R package [24]. Due to the computational complexity involved in reconstructing this quantity of

147     gene regulatory networks, the execution may take quite some time when analyzing the larger

148     networks if not distributed across a large compute cluster.

149     **RESULTS**

150     **Simulation**

151     We first tested the ENA methods on a wide array of simulated datasets. We simulated the

152     gene expression datasets based on previously observed protein-protein interaction networks[25,26],

153     and the expression data were simulated from conditional normal distributions [27].  We extract five

154     different network sizes in an approximately scale-free topology: 17 genes with 20 connections, 44

155     genes with 57 connections, 83 genes with 114 connections, 231 genes with 311 connections, or 612

156     genes with 911 connections. For each network size, we simulated datasets with differing numbers of

157     samples (microarrays): 20, 50, 100, 200, 500, and 1,000. Finally, we varied the noise by setting the

158     standard deviation of the expression values to either 0.25, 0.5, 1.0, or 1.5. In total, we generated 120

159     datasets to cover all possible arrangements of the above variables.

160     To test the effect of integrating networks derived from different datasets, we generated three

161     different datasets, each containing 200 samples, from the 231-gene networks with noise values

162     (standard deviation of the distribution of gene expression) of 0.25, 1, and 2. We then used the

163     methods described above to reconstruct three networks, one from each dataset and then aggregate

164     those networks. For comparison, we also combined all three datasets into a single dataset containing

165     these 600 samples and reconstructed a single network from this larger dataset.

166  The performance of methods in this setting can be represented on a Receiver Operating
167  Characteristic (ROC) Curve, which plots the True Positive Rate against the False Positive Rate,
168  demonstrating the performance of the method at all relevant edge weight thresholds. The
169  performance of a method can be quantified by calculating the Area Under the ROC Curve (AUC). The
170  greater the AUC, the better the performance of the method represented. A perfect reconstruction
171  would have an AUC of 1, and a random guess could obtain an AUC of 0.5.

172  **Ensemble networks derived from bootstrapping Samples**

173  We found that bootstrapping samples can increase the accuracy of network inference. For
174  example, the networks reconstructed from the dataset on the 231-gene network with a noise value of
175  0.25 can be compared to demonstrate the variations in performance as seen in Figures 1 and 2.
176  Figure 1 shows that by bootstrapping samples in the SPACE algorithm, the AUC of the reconstructed
177  network can improve from 0.75 to 0.82. Figure 2 shows the degree of AUC improvement with each
178  iteration of bootstrapping on SPACE, WGCNA and GeneNet with sample sizes of 20, 50 and 100 (left,
179  middle and right panels). From this figure, the bootstrapping method increases the performance of
180  SPACE substantially, improves GeneNet slightly when the number of microarrays is small, but does
181  not noticeably improve the performance of WGCNA. SPACE benefits from bootstrapping in 80% of all
182  simulated networks, and in 89% of "large" network simulations. Figure 3 shows the average
183  performance increase achieved by bootstrapping SPACE on different network sizes. The
184  improvement increases as the network size increases. Based on this evidence, we suggest employing
185  the bootstrapping approach when using the SPACE algorithm, but not the others evaluated in this
186  study.

187  **Ensemble networks derived from different methods**

188  Aside from optimizing individual reconstruction techniques, we find that combining different
189  network reconstruction techniques executed on the same dataset also has the power to significantly
190  improve the accuracy of the reconstructed networks. Using the dataset from the 83-gene network with
191  200 samples and a noise value of 0.25, we can review the comparative performance of each
192  reconstruction technique, as well as the aggregated network. Figure 4 shows that the aggregated
193  network outperforms any of the individual reconstruction techniques.

194  We observe that this trend holds true across most of the datasets that we tested: the
195  aggregated method typically outperforms any single reconstruction technique. This is especially
196  beneficial in scenarios in which the top performing individual network reconstruction technique may
197  vary based on the context – some methods perform well on larger networks, others excel in datasets
198  containing few samples, etc. To have an aggregation technique which consistently outperforms or
199  matches the best performing individual method eliminates the need to choose a single reconstruction
200  technique based on the context.

201  **Ensemble networks derived from different datasets**

202  Finally, we find that the ENA approach works very well when attempting to integrate various
203  datasets, especially among heterogeneous datasets that contain different distributions of expression
204  data. After generating three datasets from the 231-gene network, each with 200 samples and noise
205  values of 0.25, 1, and 2, we reconstruct each network using Bootstrapped SPACE, GeneNet, and
206  WGCNA, then aggregate the resultant networks into a single network for each dataset, producing one
207  aggregated network for each of the three datasets. We then use the ENA approach to consolidate
208  these three networks into a single network representing the underlying network behind the three
209  distinct datasets. We compare this to the alternative of simply merging all three datasets into a single
210  600-sample dataset and using the same approach to reconstruct a single network. As shown in
211  Figure 5, we find the proposed ENA approach outperforms the alternative approach of simply
212  combining the expression data into a single dataset. Reconstructing on each dataset independently
213  produces AUCs of 0.96, 0.96, and 0.89 for noise values of 0.25, 1, and 2, respectively. Naïvely
214  merging the datasets by combining them into one large dataset yields an AUC of 0.96. The network
215  aggregation approach yields the best performance, with an AUC of 0.98.

216  **Evaluating ENA approach in E. coli datasets**

217  We then tested the ENA approach on three Escherichia coli (E.coli) datasets: 1. The Many
218  Microbe Microarrays Database ("M3D")[16] contains 907 microarrays measured under 466
219  experimental conditions using Affymetrix GeneChip E.coli Genome arrays. 2. The second dataset
220  ("Str") is expression data from laboratory evolution of Escherichia coli on lactate or glycerol
221  (GSE33147)[17]. This dataset contains 96 microarrays measured under laboratory adaptive evolution
222  experiments using Affymetrix E. coli Antisense Genome Arrays. 3. The third dataset [18,19] ("BC")
223  contains 217 arrays measuring the transcriptional response of E.coli to different perturbations and
224  stresses, such as drug treatments, UV treatments and heat shock. The RegulonDB database[28,29]
225  contains the largest and best-known information on transcriptional regulation of E.coli and was used
226  as the gold standard to evaluate the accuracy of constructed networks.

227  We were able to obtain similarly positive results by employing these approaches on the E coli
228  data. Bootstrapping and aggregating the three methods on each dataset independently produced
229  AUCs of 0.574, 0.616, and 0.599 for the BC, Str, and MD3 datasets respectively. By merging the
230  three networks produced on each dataset using ENA, we were able to produce a network with an
231  AUC of 0.655, larger than the AUC of any network produced by any of the datasets independently.

232  **DISCUSSION**

233  The ability to aggregate networks using the rank-product merging approach has shown to be a
234  valuable contribution in reconstructing gene regulatory networks – and likely to other fields, as well.
235  By bootstrapping a single dataset using a single approach such as SPACE, we were able to
236  significantly improve the performance of the algorithm. By aggregating the networks produced by
237  different reconstruction techniques on a single dataset, we are able to consistently match or
238  outperform the best-performing technique for that dataset, regardless of fluctuations in the
239  performance of any one algorithm. By aggregating networks constructed independently on different

240 datasets capturing similar biological environments, we are able to reconstruct the network more
241 accurately than would be possible using any one dataset alone.

242 It is likely that SPACE was the only method to show consistent and significant improvement from
243 bootstrapping because the SPACE algorithm models the gene regulation using linear regression; as a
244 result, the network construction problem is converted to a variable selection problem. In SPACE, the
245 variable selection problem is solved by sparse regression techniques with a symmetric constraint. By
246 solving all the regression models simultaneously, SPACE is trying to get the globally optimized results.
247 However, due to the instability in variable selection [30] caused by collinearity in the data, the
248 networks constructed by SPACE are sensitive to sampling. A small change in the samples selected
249 may lead to a relatively large change in the network structure. As a result, the networks constructed
250 from bootstrapping samples are relatively "independent", which leads to better accuracy in the
251 aggregated network.

252 We provide a user-friendly R package to allow others to use these techniques on their own datasets.
253 By leveraging the MPI framework, we are able to run the bootstrapping process in parallel across
254 many cores and nodes, drastically reducing the amount of time it takes to run such analysis. We
255 include in this package a function which can permute random networks and perform ENA in order to
256 better estimate the significance of any particular connection observed in a network. This can be used
257 to reduce a continuous, complete graph to an unweighted graph including only statistically significant
258 edges.

259 Finally, we went to great lengths to ensure that all of our analysis would be as reproducible as
260 possible by structuring our analysis code in reproducible reports – most of which can be regenerated
261 at the click of a button – and making all of these freely available online at
262 https://github.com/QBRC/ENA-Research. We feel that this transparency is an important but
263 uncommon step in the scientific process and hope that other researchers begin incorporating such
264 practices in their investigation to foster more open, collaborative research.

265 **Availability**

266 The R code used to perform all of the analysis contained in this study is available in the R package
267 entitled "ENA," available on CRAN currently. The source code, as well as compiled binaries, are
268 available for download at http://cran.r-project.org/web/packages/ENA/.

269 **REFERENCES**

270 1.  Friedman N (2004) Inferring cellular networks using probabilistic graphical models. Science's STKE
271     303: 799–805.

272 2.  Stuart JM, Segal E, Koller D, Kim SK (2003) A gene-coexpression network for global discovery of
273     conserved genetic modules. Science (New York, NY) 302: 249–255. doi:10.1126/science.1087447.

274 3.  Segal E, Shapira M, Regev A, Pe D, Botstein D, et al. (n.d.) Module Networks : Discovering Regulatory
275     Modules and their Condition Specific Regulators from Gene Expression Data. Cell Research.
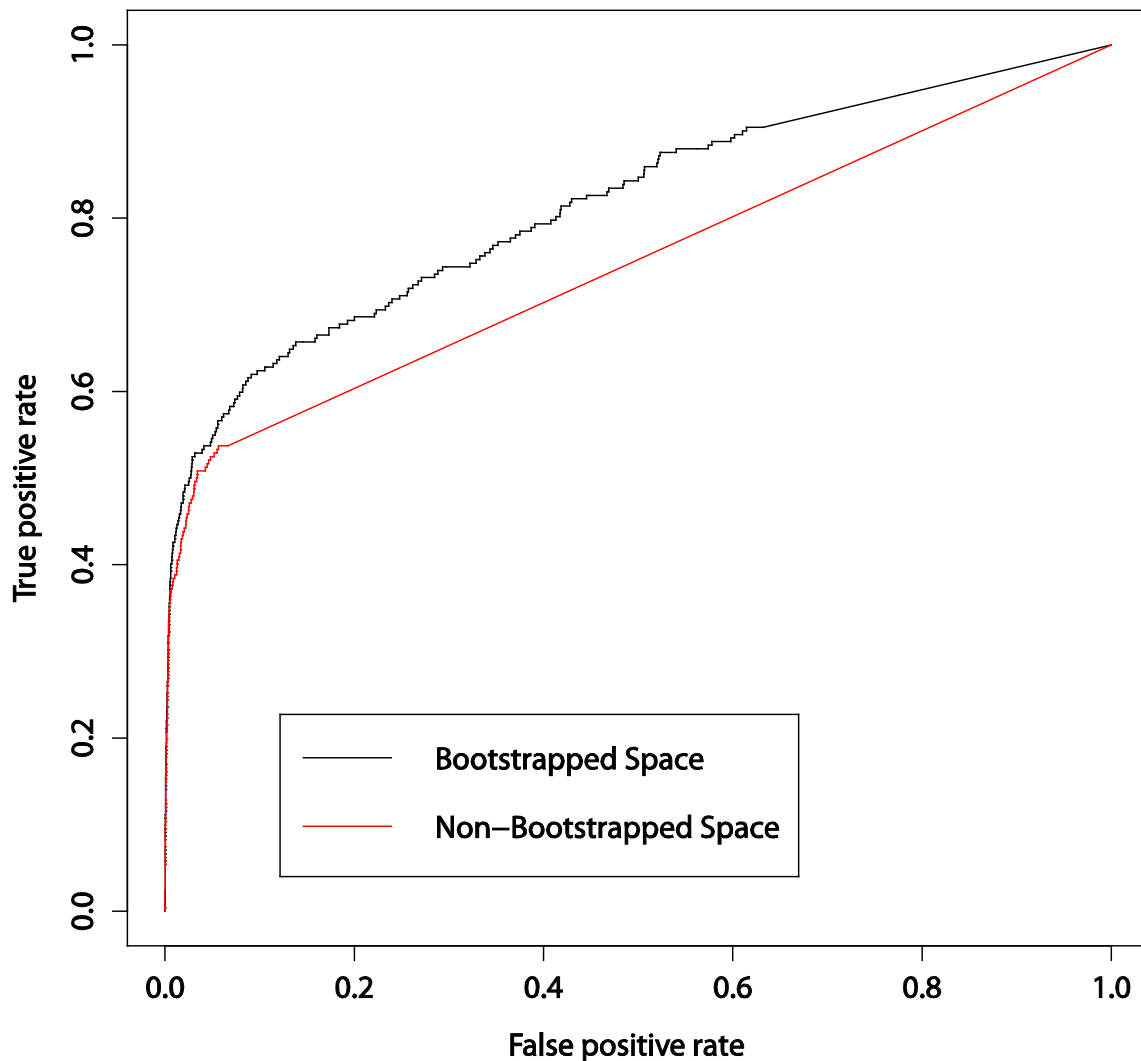
276   4.   Sachs K, Perez O, Pe'er D, Lauffenburger D a, Nolan GP (2005) Causal protein-signaling networks
277        derived from multiparameter single-cell data. Science (New York, NY) 308: 523–529.
278        doi:10.1126/science.1105809.

279   5.   Lee I, Date S V, Adai AT, Marcotte EM (2004) A Probabilistic Functional Network of Yeast Genes.
280        Science 306: 1555–1558.

281   6.   Friedman N, Linial M, Nachman I, Pe'er D (2000) Using Bayesian networks to analyze expression data.
282        Journal of computational biology : a journal of computational molecular cell biology 7: 601–620.
283        Available: http://www.ncbi.nlm.nih.gov/pubmed/11108481.

284   7.   Langfelder P, Horvath S (2008) WGCNA: an R package for weighted correlation network analysis.
285        BMC bioinformatics 9: 559. Available:
286        http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2631488&tool=pmcentrez&rendertype=abst
287        ract. Accessed 12 July 2012.

288   8.   Schäfer J, Strimmer K (2005) An empirical Bayes approach to inferring large-scale gene association
289        networks. Bioinformatics (Oxford, England) 21: 754–764. Available:
290        http://www.ncbi.nlm.nih.gov/pubmed/15479708. Accessed 12 July 2012.

291   9.   Peng J, Wang P, Zhou N, Zhu J (2007) Partial Correlation Estimation by Joint Sparse Regression
292        Model: 1–52.

293   10.  Langfelder P, Luo R, Oldham MC, Horvath S (2011) Is my network module preserved and
294        reproducible? PLoS computational biology 7: e1001057. Available:
295        http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3024255&tool=pmcentrez&rendertype=abst
296        ract. Accessed 11 June 2011.

297   11.  Breitling R, Herzyk P (2005) Rank-based methods as a non-parametric alternative of the T-statistic for
298        the analysis of biological microarray data. Journal of bioinformatics and computational …. Available:
299        http://www.worldscientific.com/doi/abs/10.1142/S0219720005001442. Accessed 22 February 2013.

300   12.  Breitling R, Armengaud P, Amtmann A, Herzyk P (2004) Rank products: a simple, yet powerful, new
301        method to detect differentially regulated genes in replicated microarray experiments. FEBS letters 573:
302        83–92. Available: http://www.ncbi.nlm.nih.gov/pubmed/15327980. Accessed 4 February 2013.

303   13.  Birmingham A, Selfors L, Forster T (2009) Statistical methods for analysis of high-throughput RNA
304        interference screens. Nature Methods 6: 569–575. Available:
305        http://www.nature.com/nmeth/journal/v6/n8/abs/nmeth.1351.html. Accessed 22 February 2013.

306   14.  Wiederhold E, Gandhi T, Permentier HP, Breitling R, Poolman B, et al. (2009) The yeast vacuolar
307        membrane proteome. Molecular & cellular proteomics : MCP 8: 380–392. Available:
308        http://www.ncbi.nlm.nih.gov/pubmed/19001347. Accessed 22 February 2013.

309   15.  Koziol J (2010) Comments on the rank product method for analyzing replicated experiments. FEBS
310        letters 584: 941–944. Available: http://www.sciencedirect.com/science/article/pii/S0014579310000542.
311        Accessed 22 February 2013.

312   16.  Faith JJ, Driscoll ME, Fusaro V a, Cosgrove EJ, Hayete B, et al. (2008) Many Microbe Microarrays
313        Database: uniformly normalized Affymetrix compendia with structured experimental metadata. Nucleic
314        acids research 36: D866–70. doi:10.1093/nar/gkm815.

315   17.  Fong SS, Joyce AR, Palsson BØ (2005) Parallel adaptive evolution cultures of Escherichia coli lead to
316        convergent growth phenotypes with different gene expression states. Genome research 15: 1365–1372.
317        doi:10.1101/gr.3832305.

318   18.  Sangurdekar DP, Srienc F, Khodursky AB (2006) A classification based framework for quantitative
319        description of large-scale microarray data. Genome biology 7: R32. doi:10.1186/gb-2006-7-4-r32.

320  19.  Xiao G, Wang X, Khodursky AB (2011) Modeling Three-Dimensional Chromosome Structures Using
321       Gene Expression Data. Journal of the American Statistical Association 106: 61–72.
322       doi:10.1198/jasa.2010.ap0950.Modeling.

323  20.  Allen JD, Xie Y, Chen M, Girard L, Xiao G (2012) Comparing Statistical Methods for Constructing
324       Large Scale Gene Networks. PLoS ONE 7: e29348. Available:
325       http://dx.plos.org/10.1371/journal.pone.0029348. Accessed 19 January 2012.

326  21.  Penrose R (1954) A Generalized Inverse for Matrices. Proc Cambridge Phil Soc 51: 406–413.

327  22.  R Core Team (2012) R: A Language and Environment for Statistical Computing.

328  23.  Gabriel E, Fagg GE, Bosilca G, Angskun T, Dongarra JJ, et al. (n.d.) Open MPI : Goals , Concept , and
329       Design of a Next Generation MPI Implementation.

330  24.  Xie Y (2012) knitr: A general-purpose package for dynamic report generation in R. Available:
331       http://cran.r-project.org/package=knitr.

332  25.  Peri S, Navarro JD, Kristiansen TZ, Amanchy R, Surendranath V, et al. (2004) Human protein reference
333       database as a discovery resource for proteomics. Nucleic acids research 32: D497–501.
334       doi:10.1093/nar/gkh070.

335  26.  Mishra GR, Suresh M, Kumaran K, Kannabiran N, Suresh S, et al. (2006) Human protein reference
336       database--2006 update. Nucleic acids research 34: D411–4. doi:10.1093/nar/gkj141.

337  27.  Pan W, Lin J, Le CT (2002) Model-based cluster analysis of microarray gene-expression data. Genome
338       biology 3: RESEARCH0009.

339  28.  Salgado H, Martı´nez-Flores I, Lo´pez-Fuentes A, Garcı´a-Sotelo JS, Liliana Porro´ n-Sotelo HS, et al.
340       (2012) Extracting Regulatory Networks of Escherichia coli from RegulonDB. In: Helden J, Toussaint A,
341       Thieffry D, editors. Bacterial Molecular Networks. New York, NY: Springer New York, Vol. 804. pp.
342       179–195. doi:10.1007/978-1-61779-361-5.

343  29.  Gama-Castro S, Salgado H, Peralta-Gil M, Santos-Zavaleta A, Muñiz-Rascado L, et al. (2011)
344       RegulonDB version 7.0: transcriptional regulation of Escherichia coli K-12 integrated within genetic
345       sensory response units (Gensor Units). Nucleic acids research 39: D98–105. doi:10.1093/nar/gkq1110.

346  30.  Breiman L (1996) Heuristics of Instability and Stabilization in Model Selection. The Annals of Statistics
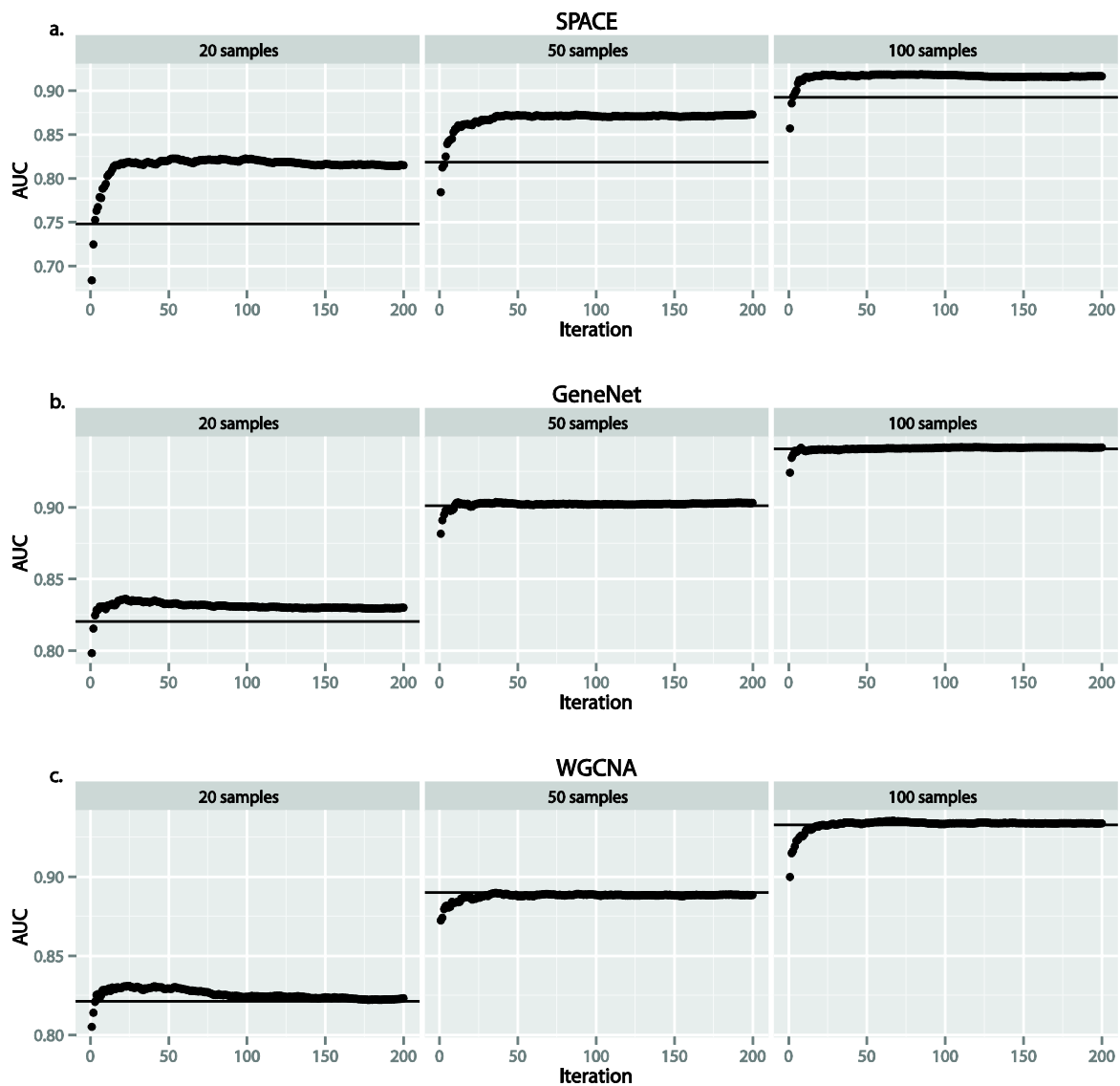347       24: 2350–2383.

348

349  **TABLE AND FIGURES LEGENDS**
350  Figure 1. Receiver Operating Characteristic (ROC) curves demonstrating the performance of the
351  SPACE algorithm on the 231-gene network with 20 samples and a noise value of 0.25 when
352  performing a single iteration or bootstrapping the dataset using the Ensemble Network Aggregation
353  approach. In this case, the Area Under the ROC Curve (AUC) of the non-bootstrapped SPACE
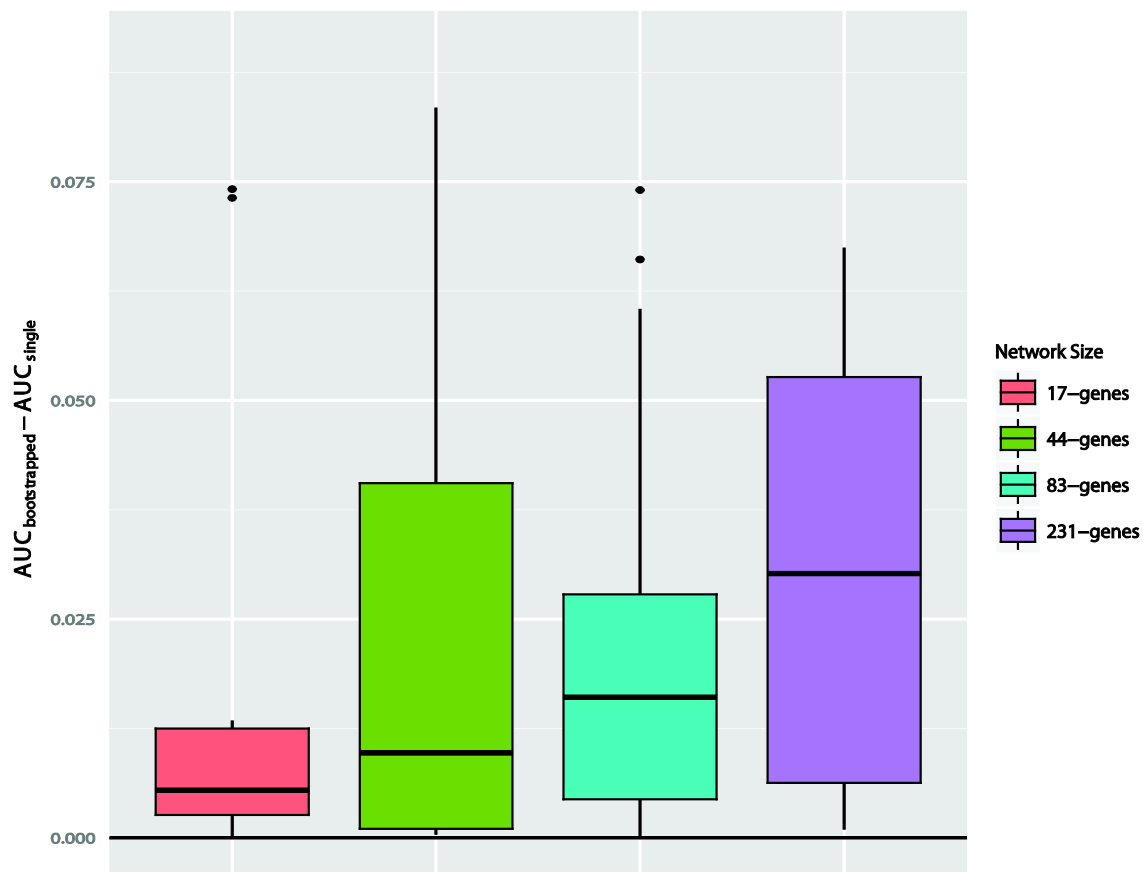354  method is 0.748, while the bootstrapped SPACE method is 0.816.

Figure 2. Comparison of the AUCs of the reconstructed networks from the 231-gene network with a noise value of 0.25 and different sample sizes (20, 50 or 100) for SPACE (a.), GeneNet (b.), and WGCNA(c.). In these plots, the y-axis shows the performance of the reconstructed network, measured by the Area Under the Curve; a horizontal line is drawn to represent the AUC of the non-bootstrapped reconstruction (a single reconstruction using all available samples). The x-axis represents the number of iterations in the bootstrapping process. Points below the horizontal line represent a loss in accuracy of the reconstructed networks, and points above the horizontal line represent a gain of AUC – an increase in performance.
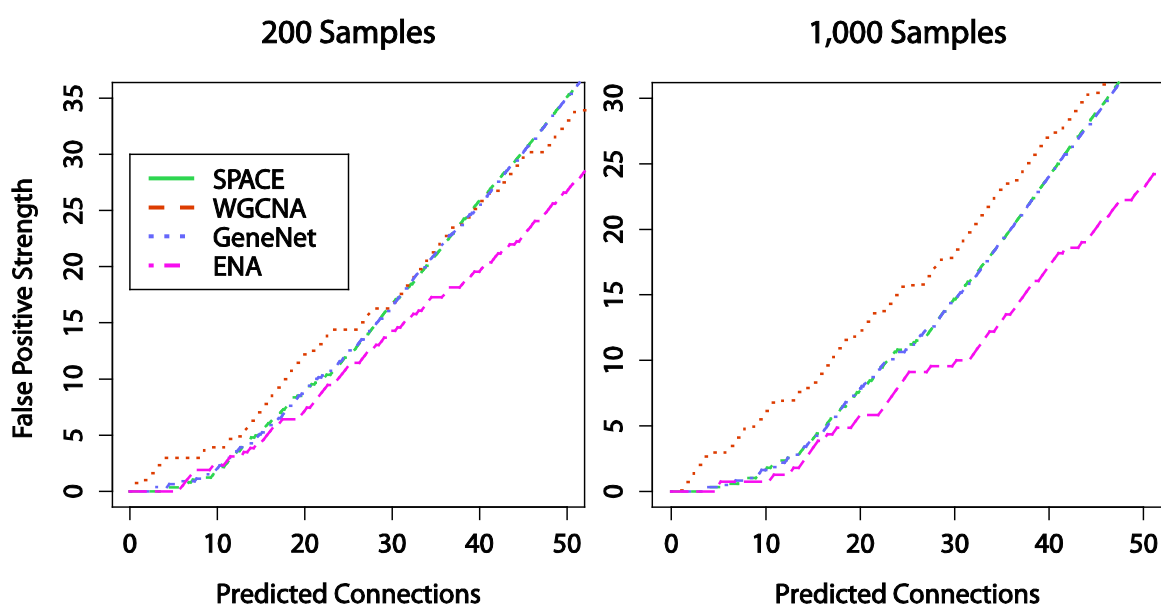
364

Figure 3. The effect of network size on ENA performance. The y-axis represents the improvement in AUC of the bootstrapped SPACE networks vs. the non-bootstrapped SPACE networks. Different bars represent different sizes of networks in the simulation study.
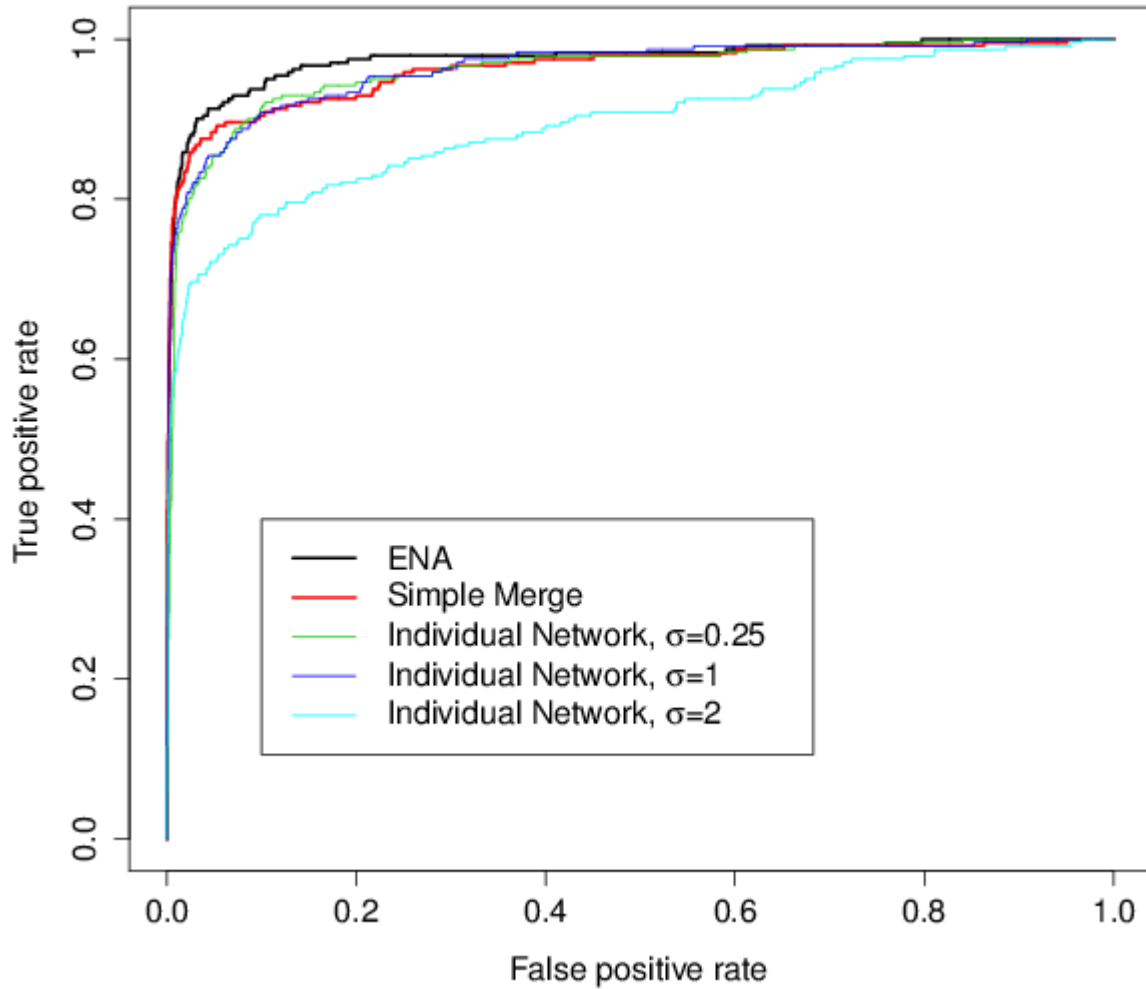
368

369 Figure 4. The performance of aggregating different methods. A comparison of the accuracy of the
370 reconstructed networks using the dataset containing 200 samples (left) and 1,000 samples (right)
371 from the 83-gene network with a noise value of 0.25. As can be seen, the ensemble network
372 aggregation approach performs better than any of the other individual techniques on these two
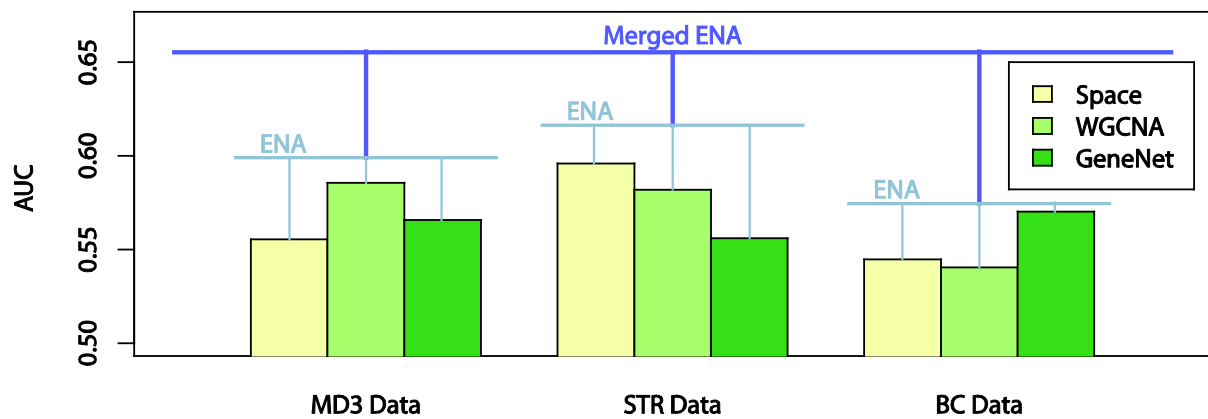373 networks.



374

375    Figure 5. The ROC curves of different approaches to reconstruct the gene network based on three
376    simulated datasets.



377

378    Figure 6. The AUCs of the produced networks when executing on the E. coli datasets. Note that the
379    aggregating networks from SPACE, WGCNA and GeneNet increases the accuracy within each
380    individual dataset, then aggregating results from three datasets further increases the accuracy
381    beyond what any one dataset offered.



382

383    Supplementary File ENA-master.zip - contains the source code for the ENA R package,

384    Supplementary File ENA-Research-Master.zip - contains all of the reproducible analysis code behind
385    this manuscript.