

1 **Supplementary Methods**

2 **General.** Unless otherwise specified, method development and analysis was performed in
3 R v3.2.5 and used Bioconductor¹.

4
5 **Pathways.** Pathway definitions were aggregated from HumanCyc² (<http://humancyc.org>),
6 IOB's NetPath³ (<http://www.netpath.org>), Reactome^{4,5} (<http://www.reactome.org>), NCI
7 Curated Pathways⁶, mSigDB⁷ (<http://software.broadinstitute.org/gsea/msigdb/>), and
8 Panther⁸ (<http://pantherdb.org/>) (downloaded from
9 http://download.baderlab.org/EM_Genesets/January_24_2016/Human/symbol/Human_Al
10 [lPathways_January_24_2016_symbol.gmt](http://download.baderlab.org/EM_Genesets/January_24_2016/Human/symbol/Human_Al))⁹. Only pathways with 10 to 500 genes were
11 included.

12
13 **Cross-method comparison.** Data consisted of TCGA level 3 gene expression data for 348
14 primary breast tumours (https://tcga-data.nci.nih.gov/docs/publications/brca_2012/);
15 the goal is binary classification of a tumour as being of type "Luminal A" or not. For all
16 methods, 67% of samples were kept as training data. The R package *caret* (v6.0-7)¹⁰ was
17 used to automate the pipeline and tune model parameters. Elastic net implementation is
18 from the R *glmnet* package (v2.0-5)¹¹, and random forests implementation from the
19 *randomForest* package (v4.6-12)¹². All three methods - elastic nets, random forests and
20 netDx - used 10-fold cross validation over three re-samplings for feature selection. Code is
21 provided as part of the netDx codebase; see [http://netdx.org/index.php/netdx-reviewer-](http://netdx.org/index.php/netdx-reviewer-page/)
22 [page/](http://netdx.org/index.php/netdx-reviewer-page/)

23
24 **Sparsification of input networks:** All negative similarities are set to zero. For a patient,
25 interactions that are not among the 50 strongest correlations, are excluded. In case of ties,
26 all interactions tied with the 50th ranked interactions are retained, for a maximum of 2% of
27 the sample size or 600 patients. This follows the parameters established in the original
28 GeneMANIA algorithm for gene expression correlation network sparsification.

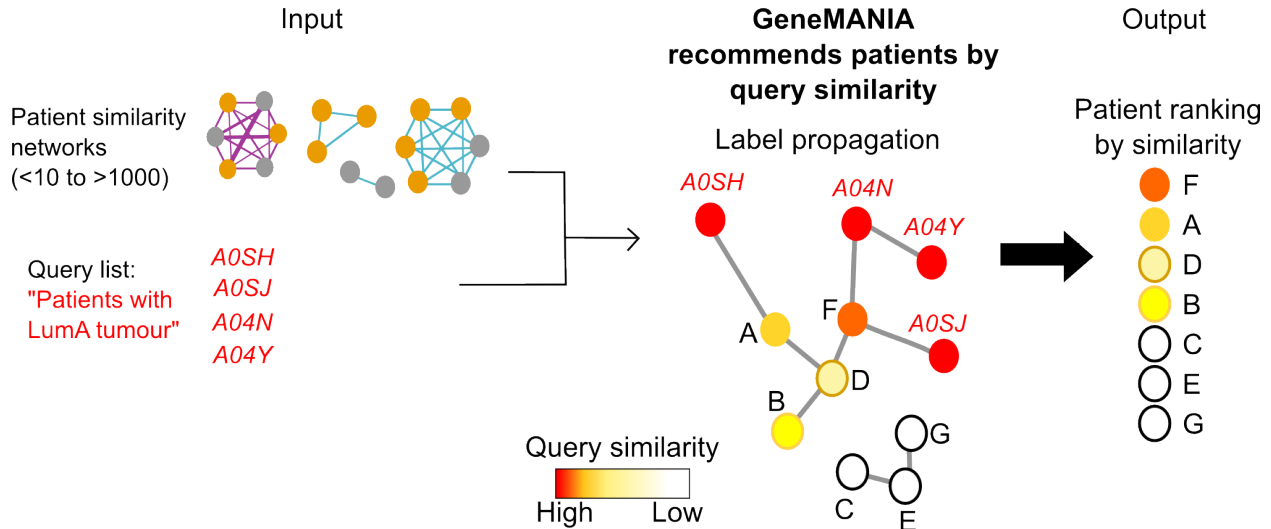
29
30 **Map of feature-selected networks.** The Enrichment Map app (v2.1.1-HOTFIX_1) in
31 Cytoscape 3.4.0 was used to generate the map of selected pathways in Figure 2D⁹. A Jaccard
32 overlap threshold of 0.05 was used to prune identical gene sets. AutoAnnotate v1.1.0 was
33 used to cluster similar pathways using MCL clustering with default parameters. The
34 network was visualized in Cytoscape 3.4.0¹³.

35
36 **Visualization of patient similarity network.** Networks that scored 10 out of 10 in the
37 feature selection algorithm were concatenated; edges with similarity less than 0.7 were
38 excluded. Multiple edges between the same pair of patients were resolved by taking the
39 maximum edge weight. The resulting network was visualized in Cytoscape. Nodes were
40 organized by the edge-weighted spring-embedded layout based on the integrated similarity
41 value.

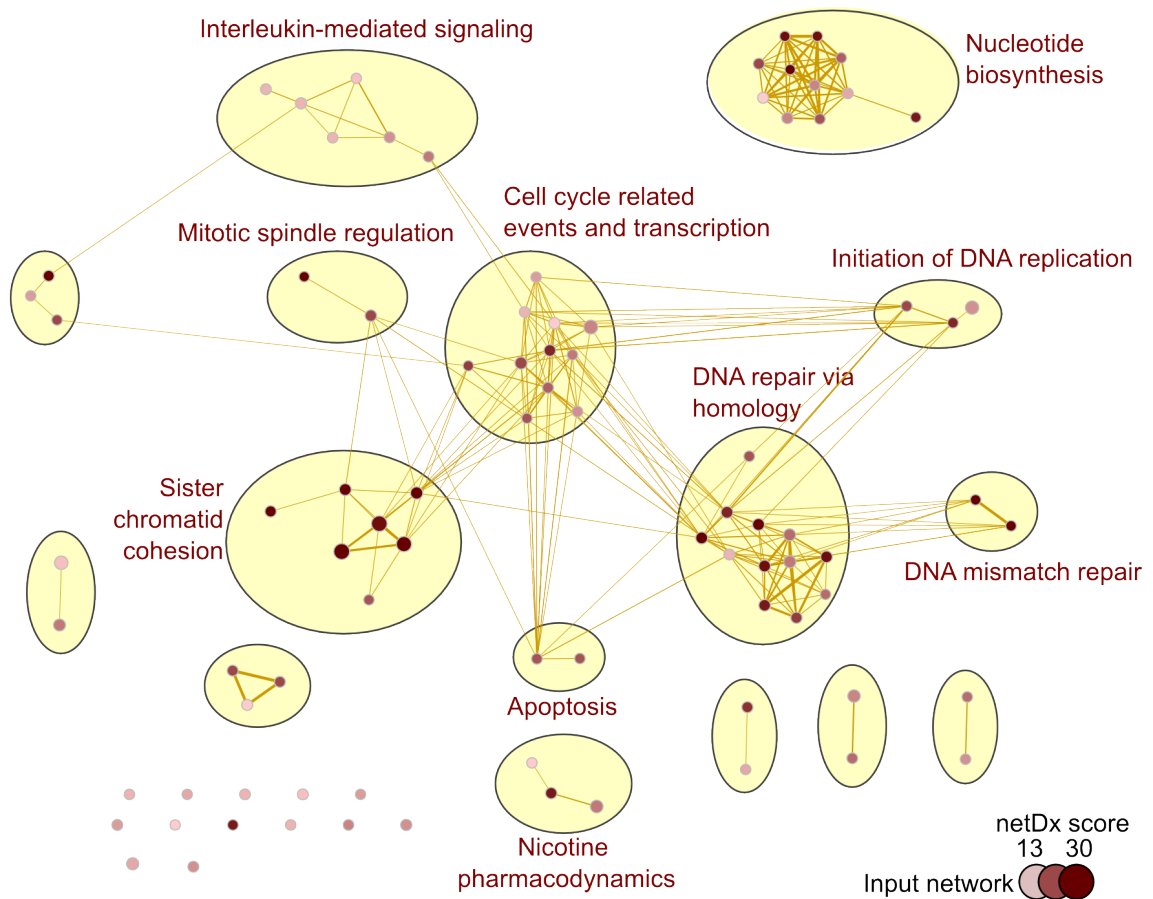
42
43 The shortest path between patient classes (a node set) was computed using Dijkstra's
44 method for weighted edges (*igraph* v1.01¹⁴); distance was defined as 1-similarity (or edge

45 weight from a patient similarity network). The overall shortest path was defined as the
46 mean pairwise shortest-path for a node set.
47
48

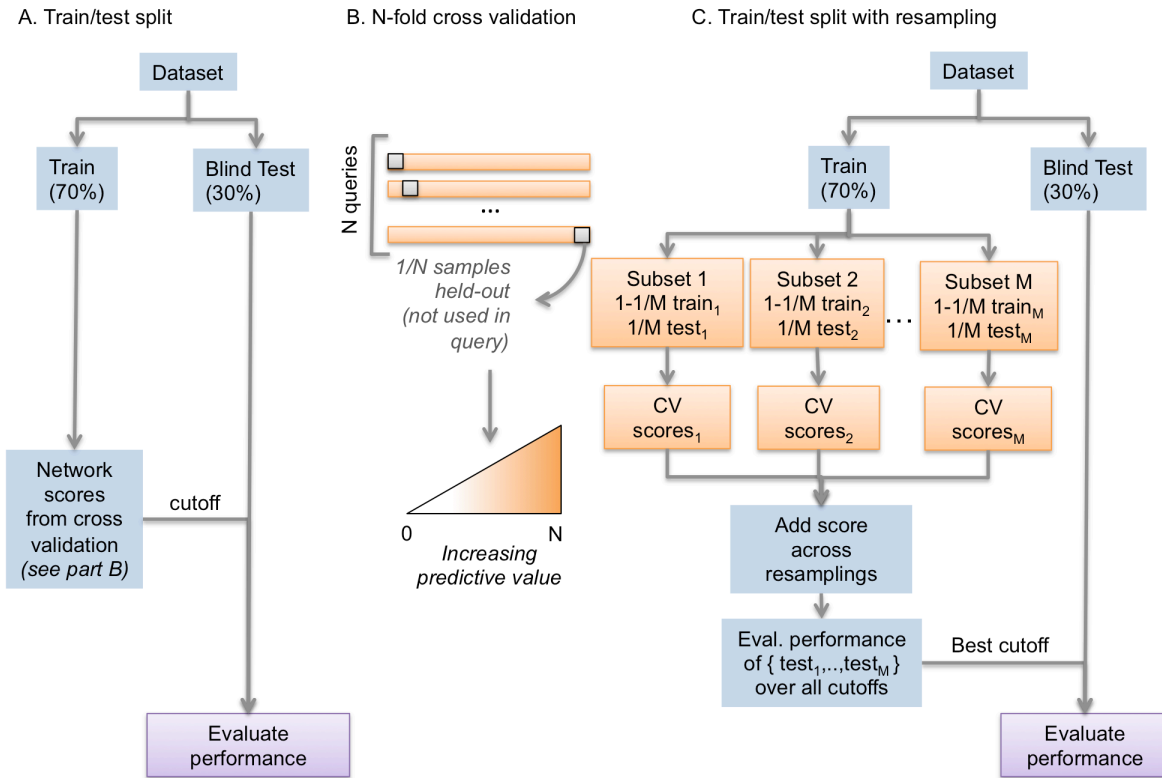
49 **Supplementary Figures**
50



51 **Supplementary Figure 1.** Conceptual overview of the GeneMANIA algorithm. GeneMANIA
52 is a network-based recommender system that ranks all nodes in its database by similarity
53 to an input query (or “positive” nodes). In the netDx application, the nodes are patients and
54 the GeneMANIA database is comprised of a set of user-defined similarity networks derived
55 from patient data (left). An example application is predicting Luminal A breast cancer type,
56 by ranking all patient tumours by similarity to known Luminal A tumours. The patient
57 ranking is achieved by a two-step process. First, input networks are integrated into a single
58 association network via regularized regression that maximizes connectivity between nodes
59 with the same label and reduces connectivity to other nodes (middle); this step computes
60 network weights or predictive value. Second, label propagation is applied to the integrated
61 network starting with the query nodes (red), thereby ranking patients from most to least
62 similar to the query (right).
63
64

66
67

68 **Supplementary Figure 2.** Networks predictive of LumA status for netDx predictor with 3-
 69 way resampling of the training data. The predictor was built using gene expression from
 70 348 primary breast tumours from the TCGA project. Each node represents a pathway
 71 (input network), and fill intensity increases with increasing netDx score. See
 72 supplementary methods for details on deriving the pathway enrichment map.
 73



74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95

Supplementary Figure 3. Variations in predictor design to reduce overfitting

- A. The data set is split into a training set and a “blind” test set. The netDx predictor is fit to the training set, and a score of predictive power is computed for each input network. Performance of the predictor is then evaluated on the test set, which is “blind” because it was not used to build the predictor. The performance is evaluated at various thresholds, to identify an optimal threshold.
- B. In the netDx algorithm, cross validation is used to score the predictive value of an input network. For a given fold size (e.g. N=10 for 10-fold cross validation), the training set is resampled N times such that each sample is excluded from exactly one set; resampling occurs without replacement. Each resampled set serves as a GeneMANIA algorithm query, and networks identified as predictive for that query have their score incremented by one. Therefore increasing N correspondingly increases the range of network scores so that finer distinctions in predictive power may be made.
- C. Another level of resampling can be introduced at the level of training samples provided for cross validation. After the initial test/train split, cross validation is computed for different random subsamples of the training set; similar to B, subsampling ensures that each sample is held out from exactly one set. The final score of a network is the cumulative score across the various iterations of feature selection. A predictor built with 3-way resampling (M=3) and 10-fold cross

96 validation results in network scores ranging from 0 to 30. Each cutoff is evaluated
97 on the test data within each resampling, and that performance is averaged. The
98 cutoff with the best mean performance is then used to validate the model on the
99 blind test.
100
101

102 **Supplementary Tables**

103

Predictor design	Input size	Runtime, workstation	Runtime, laptop	Data
Simple predictor: no resampling (network score out of 10)	348 patients in total 3,423 nets, 1-15,454 edges (mean=7,121) Of these 1,622 are sparser CNV-based nets, 1-5,671 edges (mean=186)	40 min	1h 40 min	1Gb
Complex predictor: three-way resampling (network score out of 30)	348 patients in total 1801 nets, 8,332-9,680 edges (mean=8,789)	1h 28 min	<i>Not timed</i>	1.7Gb data (2.7Mb of result files)

104

105 **Supplementary Table 1.** Computational resources required to run the breast cancer
 106 Luminal A predictor and runtimes. The runtime shown includes loading of data, building of
 107 input networks, feature selection, ranking of test samples, and performance evaluation. For
 108 laptop timings, an early 2014 MacBook Air with 1.7GHz Intel Core i7 (4 cores) with 8GB
 109 RAM was used. For workstation timings, an 2.90GHz Intel Xeon CPU (8 cores) with 128GB
 110 RAM was used.

111

112

113 **Supplementary References**

114

115 1. Gentleman, R.C. et al. Bioconductor: open software development for computational
116 biology and bioinformatics. *Genome Biol* **5**, R80 (2004).

117 2. Romero, P. et al. Computational prediction of human metabolic pathways from the
118 complete human genome. *Genome Biol* **6**, R2 (2005).

119 3. Kandasamy, K. et al. NetPath: a public resource of curated signal transduction
120 pathways. *Genome Biol* **11**, R3 (2010).

121 4. Croft, D. et al. The Reactome pathway knowledgebase. *Nucleic Acids Res* **42**, D472-7
122 (2014).

123 5. Fabregat, A. et al. The Reactome pathway Knowledgebase. *Nucleic Acids Res* **44**,
124 D481-7 (2016).

125 6. Schaefer, C.F. et al. PID: the Pathway Interaction Database. *Nucleic Acids Res* **37**,
126 D674-9 (2009).

127 7. Subramanian, A. et al. Gene set enrichment analysis: a knowledge-based approach
128 for interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A* **102**,
129 15545-50 (2005).

130 8. Mi, H. et al. The PANTHER database of protein families, subfamilies, functions and
131 pathways. *Nucleic Acids Res* **33**, D284-8 (2005).

132 9. Merico, D., Isserlin, R. & Bader, G.D. Visualizing gene-set enrichment results using
133 the Cytoscape plug-in enrichment map. *Methods Mol Biol* **781**, 257-77 (2011).

134 10. Kuhn, M. Building Predictive Models in R Using the caret Package. *Journal of*
135 *Statistical Software* **28**(2008).

136 11. Friedman J, H.T., Tibshirani R. Regularization Paths for Generalized Linear Models
137 via Coordinate Descent. *Journal of Statistical Software* **33**, 1-22 (2010).

138 12. Liaw A., W.M. Classification and Regression by randomForest. *R News* **2**, 18-22
139 (2002).

140 13. Shannon, P. et al. Cytoscape: a software environment for integrated models of
141 biomolecular interaction networks. *Genome Res* **13**, 2498-504 (2003).

142 14. Csardi G., N.T. The igraph software package for complex network research.
143 *InterJournal Complex Systems*, 1695 (2006).

144