

Supporting Information

Conaco et al. 10.1073/pnas.1201890109

SI Materials and Methods

Sponge RNA Sequencing. *Amphimedon queenslandica* samples were collected from Heron Island Reef, southern Great Barrier Reef, Queensland, Australia, and maintained by using standard protocols (1). Precompetent larvae were collected not more than 3 h after emergence from the brood chambers of the adult. Competent larvae were collected 6 h after emergence. Postlarvae exhibiting a flattened juvenile body plan were collected after 2 d of settlement on glass coverslips. Approximately 1,000 precompetent or competent larvae and 100 postlarvae were pooled to obtain sufficient RNA for sequencing. Adult material, devoid of brood chambers containing embryonic stages, was collected as a 5-mm core from apical to basal surface. RNA was extracted with TRIzol (Invitrogen) following the manufacturer's protocol. Contaminating DNA was removed using the DNasefree kit (Ambion). The poly (A) RNA fraction was enriched using the MicroPoly(A)Purist kit (Ambion) and ribosomal RNA was depleted by using the RiboMinus Eukaryote kit (Invitrogen). Poly(A) RNA fragment libraries were constructed as previously described (2) and sequenced to 50 bases by using the Applied Biosystems SOLiD V3 chemistry. Reads were mapped to assembled sponge contigs (<http://www.metazome.net/amphimedon>) by using the SOLiD RNA Analysis Pipeline Tool (<http://solidsoftwaretools.com>) and aligned using the anchor-extend method, as previously described (3). Approximately 33 to 70 million mapped reads were obtained for each developmental stage. Estimates of transcript expression were determined by counting the reads aligning to exons of predicted sponge gene models. To compare gene expression across developmental stages, read counts were normalized to the number of mapped reads per kilobase of transcript per million mapped reads per sample. Technical replicates were averaged and transcripts detected at >1 reads per kilobase of transcript per million mapped reads per sample were considered to be expressed.

Null-Model Networks. We constructed three types of null-model networks to which we compared the average correlation, R , and modularity, Q , of the true gene coexpression networks. Each of these null models was constructed to test a separate null hypothesis.

First, we constructed what we call in the main manuscript a "true random" matrix. This matrix was constructed from the following data: for each gene and each developmental stage, we drew random numbers from a uniform distribution over the interval [0,1], thereby constructing a set of N pseudogene expression patterns. We calculated the Pearson correlation coefficient between all possible pairs of pseudogene expression patterns to create a symmetric $N \times N$ matrix, $M_{\text{true random}}$, that we could directly compare with the coregulation network extracted from the true gene expression data. This process was repeated 100 times to construct 100 different $M_{\text{true random}}$ matrices, from each of which we calculated an R' value (average correlation) and Q' value (modularity). We expect the average correlation of true random networks to be small, as the correlation between pure random variables is small. Furthermore, we expect a higher modularity in true random networks (whereby multiple communities will form simply from the similarities in pseudogene expression patterns) than in highly cohesive networks in which a single underlying expression pattern dominates the matrix behavior.

Second, we constructed what we call the time-permuted control. For each gene, we randomly reassigned the expression levels to different developmental stages. This scrambling of time order creates a null model that can be used to specifically test whether the organization seen in the true data networks is dependent on

the arrow of time, from development stage 1 to n . Again, we calculated the Pearson correlation coefficient between all possible pairs of time-permuted gene expression patterns to create a symmetric $N \times N$ matrix, $M_{\text{time-permuted}}$, that we could directly compare with the coregulation network extracted from the true gene expression data. This process was repeated 1,000 times to construct 1,000 different $M_{\text{time-permuted}}$ matrices, from each of which we calculated an R' value (average correlation) and Q' value (modularity). We expect the average correlation of time-permuted networks to be small, as the correlation between time-permuted variables is small. Similar to the previous control, we further expect the modularity of time-permuted networks to be higher than that of highly cohesive networks in which a single underlying expression pattern dominates the matrix behavior.

Third, we constructed what we call the random gene set control. For each group of N genes in a given species, we extracted the expression patterns for N genes randomly selected from the entire transcriptome of that species. In effect, this creates a pseudo-gene set containing a random combination of genes expressed in a species that may or may not function together in the same cellular context. This random extraction creates a null model that can be used to specifically test whether the organization of a single functional cellular machine is different from what one would expect from a collection of genes taken from a variety of cellular machines. Again, we calculated the Pearson correlation coefficient between all possible pairs of random-gene-set gene expression patterns to create a symmetric $N \times N$ matrix, $M_{\text{random gene set}}$, that we could directly compare with the coregulation network extracted from the true gene expression data. This process was repeated 100 times to construct 100 different $M_{\text{random gene set}}$ matrices, from each of which we calculated an R' value (average correlation) and Q' value (modularity). We expect the average correlation of random gene set networks to be small, as the number of potential unique expression patterns present in the transcriptome is large. Similar to the previous two controls, we further expect the modularity of random gene set networks (in which multiple cellular machines are simultaneously probed) to be higher than that in highly cohesive networks (i.e., single functional machines) in which a single underlying expression pattern dominates the matrix behavior.

Examination of Topological Organization of Positive and Negative Correlations in Gene Expression Patterns. For each organism and each machine, we separated the coregulation matrix into a positive correlation matrix and a negative correlation matrix. The positive correlation matrix A^+ was constructed by setting all elements A_{ij}^+ equal to zero for which A_{ij} were negative and all elements A_{ij}^+ equal to A_{ij} for which A_{ij} were positive. Similarly, the negative correlation matrix A^- was constructed by setting all elements A_{ij}^- equal to zero for which A_{ij} were positive and all elements A_{ij}^- equal to $-1 * A_{ij}$ for which A_{ij} were negative. The topological organization of the positive and negative matrices (A^+ and A^- , respectively) was probed by examining the binary clustering coefficient as a function of network density. For example, for the positive correlation matrix A^+ , we constructed a binary matrix with density x by thresholding the A^+ to retain only the $x\%$ strongest (i.e., highest valued) connections. This process was performed iteratively to create a set of binary graphs in which x ranged from 0 to the maximum possible comparable density. The maximum comparable density was defined as the maximum of the following two variables: (i) the density of the graph constructed by retaining all positive correlations and (ii) the density of the graph

constructed by retaining all negative correlations. For each binary graph in this set, we computed the average clustering coefficient, C . The clustering coefficient of node i is defined by supposing that a node i has k_i neighbors, so a maximum of $k_i(k_i - 1)/2$ edges can exist between them. The local clustering coefficient C_i is the fraction of these possible edges that actually exist. The clustering coefficient of the entire network C is then defined as the mean of C_i over all nodes i . We find that the average clustering coefficient of the matrix constructed from the positive correlations, A^+ , is consistently higher than that of the matrix constructed from the negative correlations, A^- , over the full range of comparable densities (Fig. S2). This indicates that the positive correlations are more ordered than the negative correlations, and suggests that they may be more pertinent to biologically relevant phenomena.

SI Discussion

Gene Coexpression Network Analysis. Mounting evidence suggests that functional gene modules can be reliably identified by using network analysis. At least two types of network approaches are commonly applied to gene coexpression patterns: a statistical inference approach [based on Bayesian networks (4, 5)] and a data-driven exploratory approach (based on graph theory and network theory). A common type of this latter exploratory approach is weighted gene network analysis (WGNA) (6, 7), which has been used, for example, to identify functional modules in yeast (8), primate brain (9), and embryonic stem cells (10). WGNA has also been used to examine the evolution of gene coexpression modules from chimpanzees to humans, in which it facilitated the identification of modules of coexpressed genes that correspond to discrete brain regions and the quantification of their conservation between the species (9).

Like WGNA, our approach is based on representing genetic interactions as weighted networks derived from correlations in gene expression patterns. However, our approach differs from WGNA in several ways. First, we study the organization of these weighted correlation networks themselves, in which edges are based on similarities between two gene expression patterns. WGNA, on the contrary, studies the organization of a topological overlap matrix derived from the weighted correlation network, whereby edges are based on the similarities between two genes' connectivity patterns, a higher-order measure than their expression patterns. In WGNA, this topological overlap matrix is then transformed (by taking the matrix to a power b) to emphasize the importance of strong connections and deemphasize the importance of weak connections. This transformation is performed based on the assumption of an underlying scale-free topology of gene interactions, an assumption that has not yet been validated by using emerging statistical methods for the reliable identification of power laws (11, 12).

The second main difference between our approach and that of WGNA lies in the tools used to cluster the gene coexpression patterns into gene modules. WGNA employs a hierarchical clustering approach based on the dynamic branch cutting algorithm of Langfelder et al. (13), whereas we have drawn from network theory to use recently developed community detection techniques (14, 15) based on the optimization of the modularity quality function (16, 17). One of the strengths of our modularity approach lies in its ability to tune the relative importance of positive and negative correlations (using the two resolution parameters γ^+ and γ^- as discussed in the main text), and to tune the size of the modules (as demonstrated across a range of γ^+ in Fig. S3).

- Degnan BM, et al. (2009) The demosponge *Amphimedon queenslandica*: Reconstructing the ancestral metazoan genome and deciphering the origin of animal multicellularity. *Emerging Model Organisms: A Laboratory Manual* (Cold Spring Harbor Lab Press, New York), Vol 1, pp 139–166.
- Cloonan N, et al. (2008) Stem cell transcriptome profiling via massive-scale mRNA sequencing. *Nat Methods* 5:613–619.
- Tang F, et al. (2010) RNA-Seq analysis to capture the transcriptome landscape of a single cell. *Nat Protoc* 5:516–535.
- Pearl J (1988) *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference* (Morgan Kaufmann, San Francisco).
- Pearl J (2000) *Causality: Models, Reasoning, and Inference* (Cambridge Univ Press, Cambridge).
- Zhang B, Horvath S (2005) A general framework for weighted gene co-expression network analysis. *Stat Appl Genet Mol Biol* 4:Article17.
- Zhao W, et al. (2010) Weighted gene coexpression network analysis: State of the art. *J Biopharm Stat* 20:281–300.
- Carlson MR, et al. (2006) Gene connectivity, function, and sequence conservation: Predictions from modular yeast co-expression networks. *BMC Genomics* 7:40.
- Oldham MC, Horvath S, Geschwind DH (2006) Conservation and evolution of gene coexpression networks in human and chimpanzee brains. *Proc Natl Acad Sci USA* 103:17973–17978.
- Mason MJ, Fan G, Plath K, Zhou Q, Horvath S (2009) Signed weighted gene co-expression network analysis of transcriptional regulation in murine embryonic stem cells. *BMC Genomics* 10:327.
- Clauset A, Shalizi CR, Newman MEJ (2009) Power-law distributions in empirical data. *SIAM Rev* 51:661–703.
- Klaus A, Yu S, Plenz D (2011) Statistical analyses support power law distributions found in neuronal avalanches. *PLoS ONE* 6:e19779.
- Langfelder P, Zhang B, Horvath S (2008) Defining clusters from a hierarchical cluster tree: The Dynamic Tree Cut package for R. *Bioinformatics* 24:719–720.
- Fortunato S (2010) Community detection in graphs. *Phys Rep* 486:75–174.
- Porter MA, Onnela J-P, Mucha PJ (2009) Communities in networks. *Notices Am Math Soc* 56:1082–1097.
- Newman MEJ (2004) Fast algorithm for detecting community structure in networks. *Phys Rev E Stat Nonlin Soft Matter Phys* 69:066133.
- Newman MEJ, Girvan M (2004) Finding and evaluating community structure in networks. *Phys Rev E Stat Nonlin Soft Matter Phys* 69:026113.

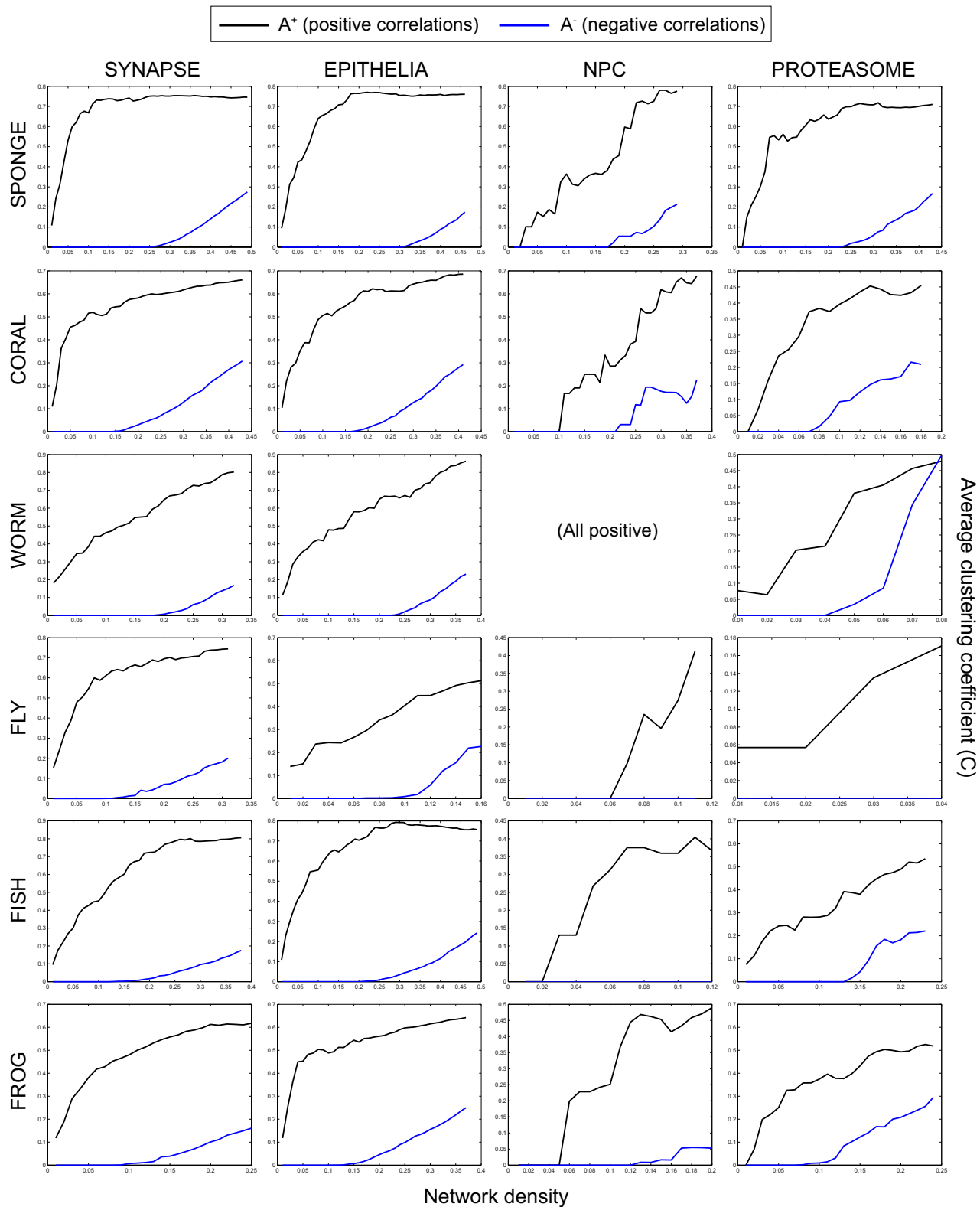


Fig. S2. Average clustering coefficient, C , as a function of the density of the gene coexpression network. The average clustering coefficient of the network matrix constructed from the positive correlations (A^+ ; black line) is consistently higher than that of the matrix constructed from the negative correlations (A^- ; blue line) over the full range of comparable densities (as defined in *Materials and Methods*).

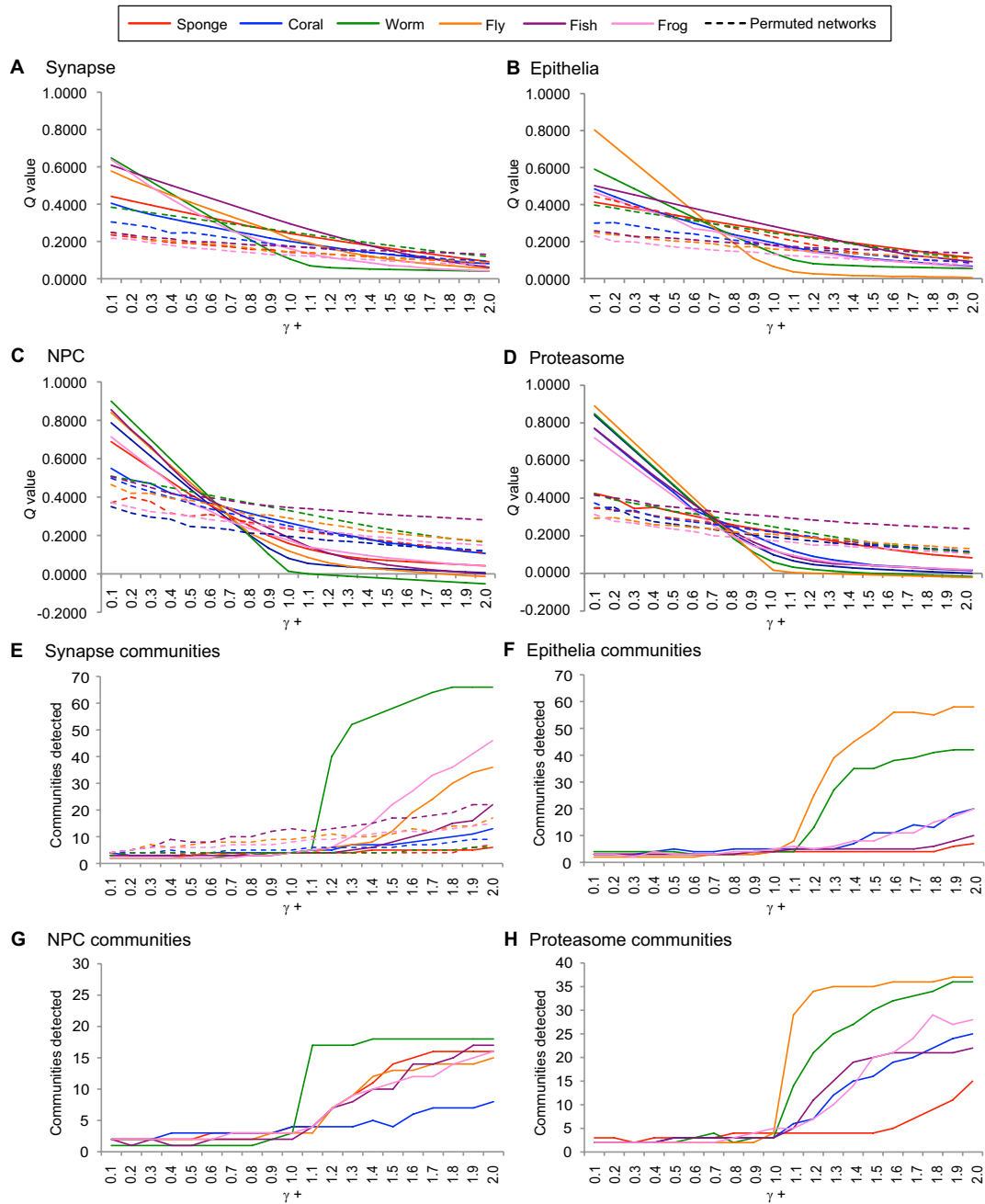


Fig. 53. Robustness of modularity estimates. Q values were computed for each network for different values of the positive resolution parameter, γ^+ , from 0.1 to 2.0 (γ^- was kept at 0.1 for the reasons outlined in the main text). Compared with synapse (A) and epithelia (B), the NPC (C) and proteasome (D) networks exhibit Q value curves (solid lines) that are more tightly distributed over species, which suggests a more cohesive structure that is also more distinct from the time-permuted networks (dashed lines). (E–H) For resolution parameters $\gamma^+ > 1.0$, Q values plateau and the number of detected communities spikes, indicating a loss of detectable community structure. Over the smaller values of the resolution parameter (e.g., $0 < \gamma^+ < 1.0$), the number of detected communities is relatively stable. Time permutation of gene expression data for synapse and epithelia resulted in networks that have, on average, lower Q values and show less sensitivity to changes in γ^+ .

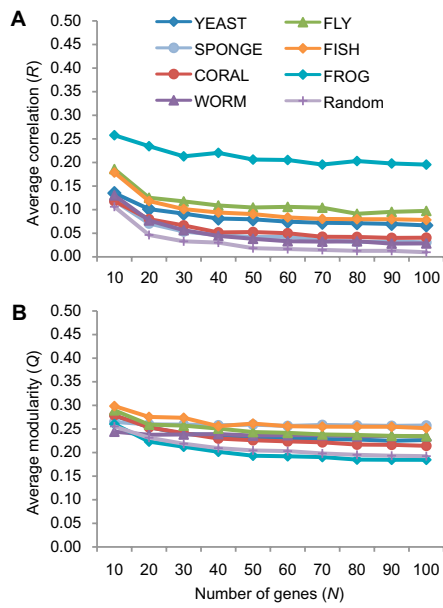


Fig. S4. Average correlation (A) and modularity (B) of 100 sets of N genes randomly selected from the transcriptome datasets of six organisms or size-matched random number matrices (random). Random number matrices tend to display lower R and Q in comparison with the matrices derived from random sets of genes.

Table S3. Results for the modularity analysis using an alternative heuristic (spectral optimization method)

Species	<i>N</i>	Modularity, $Q_{spec} \pm SD^*$	Communities (<i>S</i>)	$Q_{spec_rand} \pm SD^\dagger$	<i>P</i> value (Q_{spec} vs. Q_{spec_rand}) [‡]
Synapse					
Sponge	87	0.235 ± 3.93E-16	4	0.238 ± 0.016	0.82
Coral	87	0.186 ± 3.04E-16	4	0.200 ± 0.014	0.34
Worm	107	0.101 ± 2.15E-16	4	0.224 ± 0.018	6.28E-10
Fly	112	0.222 ± 4.50E-16	4	0.220 ± 0.017	0.91
Fish	84	0.279 ± 4.87E-16	3	0.237 ± 0.022	0.07
Frog	107	0.149 ± 2.31E-16	4	0.171 ± 0.015	0.14
Epithelia					
Sponge	70	0.248 ± 3.63E-16	4	0.239 ± 0.018	0.62
Coral	74	0.164 ± 2.04E-16	3	0.202 ± 0.018	0.0400
Worm	68	0.132 ± 2.05E-16	5	0.224 ± 0.020	1.81E-05
Fly	75	0.064 ± 1.10E-16	4	0.224 ± 0.020	1.19E-12
Fish	76	0.272 ± 4.48E-16	4	0.244 ± 0.021	0.20
Frog	99	0.168 ± 2.62E-16	4	0.172 ± 0.012	0.74
NPC					
Yeast	30	0.074 ± 1.13E-16	5	0.229 ± 0.026	3.31E-08
Sponge	23	0.138 ± 2.34E-16	3	0.253 ± 0.029	1.64E-04
Coral	14	0.267 ± 5.53E-16	4	0.260 ± 0.028	0.82
Worm	18	0.013 ± 6.07E-17	3	0.389 ± 0.086	3.04E-05
Fly	17	0.120 ± 1.43E-16	3	0.263 ± 0.035	8.82E-05
Fish	23	0.190 ± 9.55E-17	2	0.258 ± 0.034	0.05
Frog	26	0.167 ± 2.47E-16	4	0.211 ± 0.027	0.12
Proteasome					
Yeast	37	0.092 ± 9.93E-17	2	0.226 ± 0.024	2.62E-07
Sponge	39	0.204 ± 1.89E-16	5	0.244 ± 0.025	0.12
Coral	38	0.156 ± 2.36E-16	3	0.223 ± 0.019	7.29E-04
Worm	39	0.060 ± 8.94E-17	3	0.232 ± 0.025	7.61E-10
Fly	38	0.016 ± 6.53E-17	3	0.238 ± 0.023	5.98E-16
Fish	31	0.118 ± 2.01E-16	3	0.257 ± 0.027	1.99E-06
Frog	46	0.114 ± 1.79E-16	5	0.189 ± 0.020	2.15E-04

The alternative heuristic (spectral optimization) is a community detection method developed by Newman (1), 2006 ($\gamma^+ 1.0$, $\gamma^- 0.1$).

**Q* values computed using the Blondel et al. (2), 2008 vs. spectral optimization methods are highly correlated ($r = 0.9960$, $P = 1.17E-26$).

[†]Computed from 100 sets of *N* genes randomly selected from the whole transcriptome.

[‡]Two-tailed *t* test.

1. Newman MEJ (2006) Finding community structure in networks using the eigenvectors of matrices. *Phys Rev E Stat Nonlin Soft Matter Phys* 74:036104.
2. Blondel VD, Guillaume J-L, Lambiotte R, Lefebvre E (2008) Fast unfolding of communities in large networks. *J Stat Mech* 10:P10008.

Other Supporting Information Files

[Dataset S1 \(XLS\)](#)