

# The topological relationship between the large-scale attributes and local interaction patterns of complex networks

A. Vázquez\*, R. Dobrin†, D. Sergi‡, J.-P. Eckmann\*§, Z. N. Oltvai†, and A.-L. Barabási\*¶

\*Department of Physics and Center for Complex Network Research, University of Notre Dame, Notre Dame, IN 46556; †Department of Pathology, Northwestern University, Chicago, IL 60611; and ‡Département de Physique Théorique and §Section de Mathématiques, Université de Genève, CH-1211 Geneva, Switzerland

Edited by Harry L. Swinney, University of Texas, Austin, TX, and approved November 1, 2004 (received for review August 16, 2004)

Recent evidence indicates that the abundance of recurring elementary interaction patterns in complex networks, often called subgraphs or motifs, carry significant information about their function and overall organization. Yet, the underlying reasons for the variable quantity of different subgraph types, their propensity to form clusters, and their relationship with the networks' global organization remain poorly understood. Here we show that a network's large-scale topological organization and its local subgraph structure mutually define and predict each other, as confirmed by direct measurements in five well studied cellular networks. We also demonstrate the inherent existence of two distinct classes of subgraphs, and show that, in contrast to the low-density type II subgraphs, the highly abundant type I subgraphs cannot exist in isolation but must naturally aggregate into subgraph clusters. The identified topological framework may have important implications for our understanding of the origin and function of subgraphs in all complex networks.

aggregation | subgraphs

A number of complex biological and nonbiological networks were recently found to contain network motifs, representing elementary interaction patterns between small groups of nodes (subgraphs) that occur substantially more often than would be expected in a random network of similar size and connectivity (1, 2). Theoretical and experimental evidence indicates that at least some of these recurring elementary interaction patterns carry significant information about the given network's function and overall organization (1–4). For example, transcriptional regulatory networks of cells (1, 2, 5, 6), neural networks of *C. elegans* (2), and some electronic circuits (2) are all information processing networks that contain a significant number of feed-forward loop (FFL) motifs. However, in transcriptional regulatory networks these motifs do not exist in isolation but meld into motif clusters (7), while other networks are devoid of FFLs altogether (2).

In general, all subgraphs have two important properties: their topology and the directionality of their links. In cellular networks, these two properties can be clearly separated from each other. In protein–protein interaction (PPI) networks all links are by definition nondirectional. In contrast, in transcriptional regulatory networks information flow between a transcription factor and the operon (gene) regulated by it is almost always unidirectional (1, 2). Metabolic networks occupy an intermediate position between these two extremes, because most, but not all, metabolic reactions are reversible under various growth conditions. Despite the difference in the relative role of link directionality, the large-scale organization of the three different network types is quite similar, most being characterized by a scale-free connectivity distribution and hierarchical modularity (8–12). The only exception is the incoming degree distribution (i.e., the number of transcription factors regulating a target gene) of regulatory networks, which decays faster than a power law, because the number of transcription factors that can

simultaneously bind to a target gene's promoter region appears to be limited by structural constraints (13).

A coherent understanding of a network's topological and functional organization requires the development of a single framework that can explain the appearance of subgraphs and motifs, the mechanisms responsible for their aggregation into larger superstructures, and their relationship with the universal large-scale features of complex networks. Here we present such a unifying framework by focusing on five well characterized cellular networks of a prokaryotic model organism and a eukaryotic model organism, the metabolic and transcriptional regulatory networks of *Saccharomyces cerevisiae* and *Escherichia coli*, respectively, and the PPI network of *S. cerevisiae*. We show that the subgraph density in these networks can be fully predicted based on knowledge of the two parameters characterizing their global scale-free and hierarchical topology. Furthermore, we demonstrate that a network's large-scale topological organization and its local subgraph structure mutually define and predict each other. We also show the inherent existence of two distinct classes of subgraphs, demonstrating that in contrast to the low-density type II subgraphs, the highly abundant type I subgraphs cannot exist in isolation but must naturally aggregate into subgraph clusters. These results imply a fundamental unity in the origin of subgraphs and subgraph clusters in all complex networks.

## Materials and Methods

The transcriptional regulatory networks of *E. coli* and *S. cerevisiae* (1, 2) are available from [www.weizmann.ac.il/mcb/UriAlon](http://www.weizmann.ac.il/mcb/UriAlon). We have studied their undirected representations, where transcription factors and genes are represented by nodes and each regulation-based interaction is replaced by an undirected link. The metabolic networks of *E. coli* and *S. cerevisiae* were obtained from the WIT/ERGO database (14) (<http://igweb.integratedgenomics.com/IGwit>). Metabolites are represented by nodes, and undirected links connect each substrate to each product of the same reaction. The PPI network of *S. cerevisiae* was obtained from DIP (15) (<http://dip.doe-mbi.ucla.edu>). Proteins are represented by nodes, and each pairwise protein interaction is represented by an undirected link.

## Results







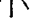

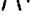

**The Abundance of Subgraphs in Cellular Networks.** Table 1 lists the density of several  $n$ -node subgraphs of the five studied intracellular molecular interaction networks: the metabolic and transcriptional regulatory networks of *S. cerevisiae* and *E. coli* and the PPI network of *S. cerevisiae*. Our study is limited to subgraphs with  $n$  nodes and  $m$  links that can be decomposed into a central node with  $n - 1$  neighbors, the remaining  $m - n + 1$  links connecting these

This paper was submitted directly (Track II) to the PNAS office.

Abbreviations: FFL, feed-forward loop; PPI, protein–protein interaction.

¶To whom correspondence should be addressed. E-mail: [alb@nd.edu](mailto:alb@nd.edu).

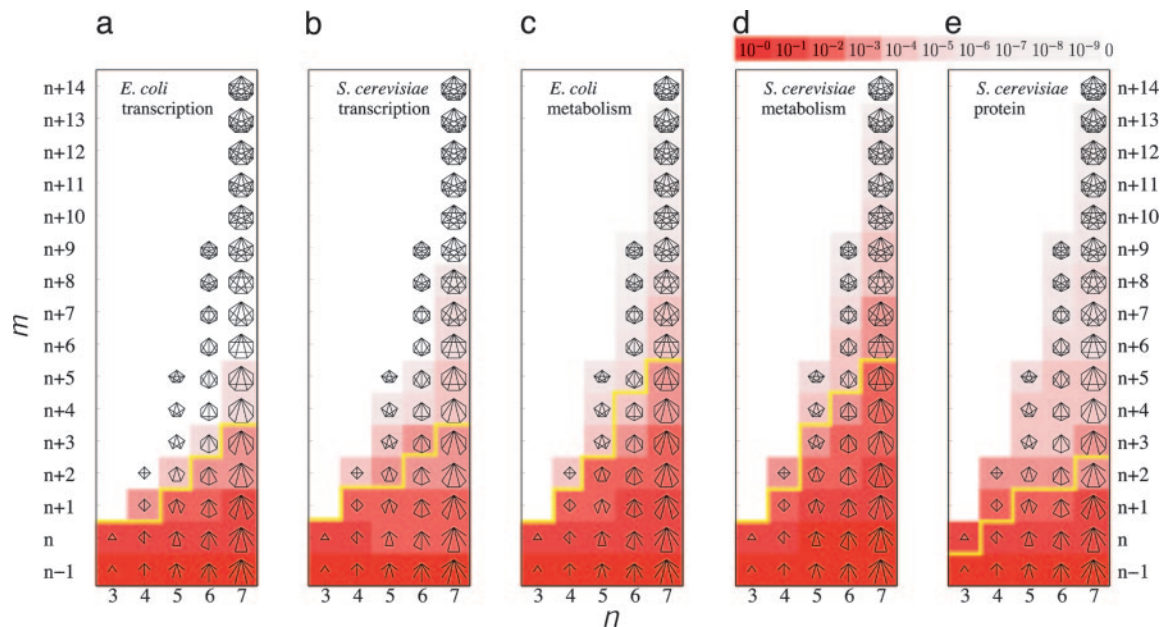
© 2004 by The National Academy of Sciences of the USA

$(n, m)$		Transcription		Metabolic		Protein Interaction
		<i>E. coli</i>	<i>S. cerevisiae</i>	<i>E. coli</i>	<i>S. cerevisiae</i>	<i>S. cerevisiae</i>
(3,2)		12	19	101	72	70
(3,3)		0.30	0.31	5.0	5.8	4.1
(4,3)		169	220	4,412	2,041	2,395
(4,6)		0.00	0.00	0.44	0.77	0.97
(5,4)		2,492	2,587	$2.1 \times 10^5$	$5.9 \times 10^4$	$1.2 \times 10^5$
(5,10)		0.00	0.00	0.055	0.20	0.66
(6,5)		$3.2 \times 10^4$	$2.8 \times 10^4$	$8.8 \times 10^6$	$1.5 \times 10^6$	$5.7 \times 10^6$
(6,15)		0.00	0.00	0.00	0.03	0.36
(7,6)		$3.4 \times 10^5$	$2.7 \times 10^5$	$3.5 \times 10^8$	$3.7 \times 10^7$	$2.4 \times 10^8$
(7,21)		0.00	0.00	0.00	0.00	0.00

exponent ( $\gamma$ ) characterizes the number of interactions in which a node is engaged, capturing the overall inhomogeneity in the connectivity of complex cellular networks: Whereas most molecules are engaged in only a few interactions, a few hubs are linked to a significantly higher number of other molecules (nodes). These wide degree variations are captured by the degree distribution, which for the studied cellular networks follows a power law,  $P(k) \sim k^{-\gamma}$  (7, 13, 20–23). In contrast, the hierarchical exponent ( $\alpha$ ) characterizes the networks’ innate modularity, indicating that many small, highly interconnected groups of nodes form larger but less cohesive topological modules (7, 19). This hierarchical modularity is captured by the scaling law (24, 38)  $C(k) \sim C_0 k^{-\alpha}$ , where  $C(k) = 2T(k)/k(k-1)$  is the clustering coefficient of a node with  $k$  links, denoting the probability that a node’s neighbors are linked to each other (25), and  $T(k)$  is the number of direct links between the node’s  $k$  neighbors. Empirical studies indicate that each cellular network is characterized by a unique pair of ( $\gamma, \alpha$ ) parameters, listed in Table 2, which were determined from the scaling of  $P(k)$  and  $C(k)$

Exponent	Transcription		Metabolic		Protein Interaction
	<i>E. coli</i>	<i>S. cerevisiae</i>	<i>E. coli</i>	<i>S. cerevisiae</i>	<i>S. cerevisiae</i>
$\gamma$	$2.1 \pm 0.3$	$2.0 \pm 0.2$	$2.0 \pm 0.4$	$2.0 \pm 0.1$	$2.4 \pm 0.4$
$\alpha$	$1.0 \pm 0.2$	$1.0 \pm 0.2$	$0.8 \pm 0.3$	$0.7 \pm 0.3$	$1.3 \pm 0.5$
$\beta$					
Meas.	$1.0 \pm 0.2$	$0.8 \pm 0.2$	$1.1 \pm 0.2$	$1.4 \pm 0.2$	$0.7 \pm 0.2$
Pred.	0.97	0.95	1.2	1.3	0.7
$\delta$					
Meas.	$2.1 \pm 0.2$	$2.2 \pm 0.2$	$1.8 \pm 0.2$	$1.7 \pm 0.2$	$2.3 \pm 0.2$
Pred.	2.0	1.9	1.8	1.8	3.0

PNAS | December 28, 2004 | vol. 101 | no. 52 | 17941



**Fig. 1.** Subgraph phase diagrams. The phase diagrams organize the subgraphs based on the number of nodes ( $n$ , horizontal axis) and the number of links ( $m$ , vertical axis), each discrete point explicitly depicting the corresponding subgraph. The stepped yellow line corresponds to the predicted phase boundary separating the abundant type I subgraphs (below the line) from the constant density type II subgraphs (above the line). The background color is proportional to the relative subgraph count  $C_{nm} = N_{nm}/\sum_s N_{ns}$  of each  $n$ -node subgraph, the color code being shown in the upper right corner. Note that some  $(n, m)$  points in the phase diagram may correspond to several topologically distinguishable subgraphs. For simplicity, we depict only one representative topology in such cases. Because the yellow phase boundary depends on the  $\gamma$  and  $\alpha$  exponents of the corresponding network, each phase diagram is slightly different. Yet, there is a visible similarity between the networks of the same kind: The phase diagrams of the two transcription or the two metabolic networks are almost indistinguishable.

functions describing the undirected version of these networks (7, 19).

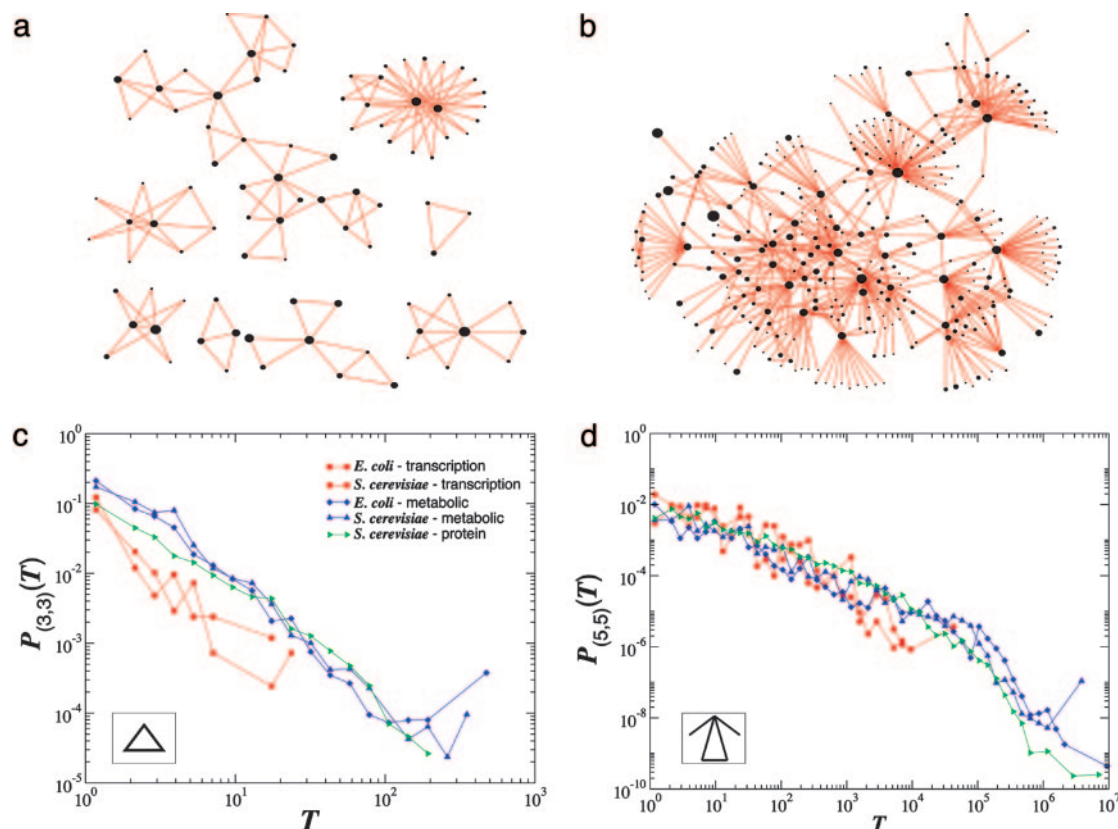
**Type I and II Subgraphs.** To examine the relationship between these two parameters and the observed subgraph density, we calculated analytically the number  $N_{nm}$  of subgraphs with  $n$  nodes and  $m$  interactions expected for a network of  $N$  nodes, in which the nodes, apart from fixed  $(\gamma, \alpha)$  parameters, are randomly connected to each other. As each pair of neighbors of a node with degree  $k$  is connected with a probability  $C(k) \sim k^{-\alpha}$ , the average number of  $(n, m)$  subgraphs that pass by a node with degree  $k$  scales as  $N_{nm}(k) \sim k^{n-1-(m-n+1)\alpha}$ . Summing over the degree distribution, we obtain the number of  $(n, m)$  subgraphs,  $N_{nm} \sim N \sum_k P(k) N_{nm}(k)$ . The convergence of this sum predicts the existence of two subgraph classes. Type I subgraphs are those that satisfy  $(m - n + 1)\alpha - (n - \gamma) < 0$ , their number being given by  $N_{nm}^I \sim N k_{\max}^{-(m-n+1)\alpha - (n-\gamma)}$ , where  $k_{\max}$  denotes the degree of the most connected node in the network. Type II subgraphs are those that satisfy  $(m - n + 1)\alpha - (n - \gamma) > 0$ , and their number is given by  $N_{nm}^{II} \sim N$ . As even for finite networks  $k_{\max} \gg 1$ , the typical number of type I subgraphs is significantly larger than the number of type II subgraphs ( $N_{nm}^I/N_{nm}^{II} \gg 1$ ). Moreover, for infinite systems ( $N \rightarrow \infty$ ) the relative number of type II subgraphs is vanishingly small compared with type I subgraphs, as  $N_{nm}^I/N_{nm}^{II} \rightarrow \infty$ . Table 1 supports these predictions, indicating that the density of the subgraphs with a minimal number of connections (extreme type I) (4,3), (5,4), (6,5), (7,6) is in the range 10 to  $10^5$  ( $N_{nm}^I \gg 1$ ). In contrast, the density of the subgraphs with a maximal number of connections (extreme type II) (4,6), (5,10), (6,15), (7,21) is either zero or close to zero, and always negligible compared with their type I counterparts.

The main results of our analysis are summarized in the  $(n, m)$  phase diagrams of Fig. 1, in which each square corresponds to a different subgraph. The  $(m - n + 1)\alpha - (n - \gamma) = 0$  condition, predicted to separate the type I and II subgraphs, appears as stepped yellow phase boundaries in the phase diagrams. For

example, for the *E. coli* transcriptional regulatory network with  $\alpha = 1$  and  $\gamma = 2.1$  (Table 2) the phase boundary corresponds to a stepped-line with approximate overall slope  $1 + 1/\alpha = 2.0$  and intercept  $-1 - \gamma/\alpha = -3.1$  (Fig. 1a). The type II subgraphs are those above this boundary and should be either absent or present only in very low numbers in the transcriptional regulatory network. In contrast, the type I subgraphs below the boundary are predicted to be abundant.

To visually highlight the validity of these predictions, we color-coded Fig. 1 according to the normalized count of each subgraph in each cellular network. We find a good agreement between the analytical predictions and the measured subgraph count: The normalized count of the type I subgraphs below the phase boundary is in the  $10^{-2}$  to 1 range, in contrast with the type II subgraphs above the predicted boundary, whose normalized count is either zero or in the  $10^{-9}$  to  $10^{-3}$  range. Comparing Fig. 1 a–e indicates that whereas the stepped phase boundaries for the different cellular networks differ because of the differences in the  $(\gamma, \alpha)$  exponents (Table 2), the observed densities in the real networks follow relatively closely the predicted phase boundaries. Occasional local deviations from the predictions can be attributed to the error bars of the  $(\gamma, \alpha)$  exponents (Table 2), which allow for some local uncertainties for the phase boundary. Fig. 1 a–e also indicates that, in agreement with the empirical findings (1–4), each cellular network is characterized by a distinct set of overrepresented type I subgraphs, raising the possibility of classifying networks based on their local structure (4). Yet, the phase diagrams demonstrate that knowledge of two global topological parameters automatically uncovers the local structure of cellular networks, suggesting that a subgraph- or motif-based classification could be equivalent with a classification based on the different  $(\gamma, \alpha)$  exponents characterizing these networks.

**Subgraphs and Motifs.** The concept of *motifs* was recently introduced to denote those subgraphs whose number exceeds by a preset



**Fig. 2.** Subgraph distributions in cellular networks. *a* and *b* show all nodes in the *S. cerevisiae* transcription regulatory network that participate in triangle (3,3) and (5,5) subgraphs (depicted in *Insets* of *c* and *d*). The size (area) of each node is drawn proportional to its degree,  $k$ , in the full network, indicating that subgraphs tend to aggregate around the hubs. Indeed, although there are hubs that have only a few subgraphs around them, in most cases subgraph aggregation is seen only around highly connected nodes. Note that the (3,3) subgraphs of *S. cerevisiae* is above the percolation boundary (Fig. 1*e*), and therefore they are broken into small islands. In contrast, the (5,5) subgraph is well below the boundary, forming a fully connected giant component, with no isolated subgraphs, as predicted. *c* and *d* show the  $P(T)$  distribution of the number of (3,3) and (5,5) subgraphs, respectively, passing by a given node. The plot indicates that for both subgraphs,  $P(T)$  approximates a power law  $P(T) \sim T^{-\delta}$ . Note the quite extended scaling regimes for some networks; e.g., for the (5,5) subgraph the scaling extends over four to five orders of magnitude. The  $\delta$  exponents measured and predicted for each network are summarized in Table 2 and the supporting information.

threshold their expected count in a randomized network (1–4). Our results indicate that overrepresented type I subgraphs are innate topological features of complex networks, and we do not need to invoke a comparison to a randomized graph or introduce a threshold parameter to identify them. Indeed, the signature of type I subgraphs is that their density increases with the number of nodes in the network ( $N_{nm}^I/N \rightarrow \infty$  as  $N \rightarrow \infty$ ), compared with the type II subgraphs, whose density is independent of the network size ( $N_{nm}^{II}/N \rightarrow \text{const}$ ). The existence of the type II subgraphs is intertwined with the network's global hierarchical topology: The decreasing  $C(k)$  reduces the likelihood that the neighbors of a highly connected node are linked to each other, therefore limiting the chance that these nodes participate in highly connected subgraphs. If  $C(k)$  were independent of  $k$  (i.e.,  $\alpha = 0$ ), only type I subgraphs would exist, since in the  $\alpha \rightarrow 0$  limit the  $1 + 1/\alpha$  slope of the yellow phase boundary diverges, eliminating all type II subgraphs. Because the absolute count of the subgraphs is the most fundamental quantity for evaluating a local interaction pattern's topological role in a network, we will continue focusing on the direct subgraph count, limiting the discussion on motifs and the role of the randomized reference frame to the supporting information, which is published on the PNAS web site. Note that the scaling of the subgraph density with the network size  $N$  was already predicted in ref. 26. Yet, the calculation did not take into account the scaling of the clustering coefficient; thus, the results are limited to the  $\alpha = 0$  limit of our predictions. Thanks to the  $C(k)$  scaling, however, for

realistic  $\gamma$  values we predict a new phase, which contains the type II subgraphs.

**Subgraphs Aggregate Around Hubs.** The very large densities we observe for some type I subgraphs (Tables 1 and 2) require us to explain how to distribute as many as  $10^{11}$  subgraphs in a network with only  $10^3$  nodes. We address this question by calculating the number of distinct subgraphs in which a given node (gene, metabolite, or protein) participates. We first focus on the triangle subgraph (3,3), the elementary building block of many higher-order subgraphs. A node with  $k$  links participates on average in  $T(k) = C(k)k(k-1)/2$  triangles. For large  $k$  this scales as  $T(k) \sim k^{2-\alpha}$ . Therefore, the probability that exactly  $T$  triangles pass through a node is  $P(T) \sim T^{-\delta}$ , where  $\delta = 1 + (\gamma - 1)/(2 - \alpha)$ , a power-law dependence that indicates that whereas the majority of nodes participate in at most one or two triangles, a few nodes take part in a very large number of triangle subgraphs. The monotonic nature of  $T(k)$  indicates that the triangles are not distributed uniformly within the network but tend to aggregate around the hubs. Because a node with  $k$  links can carry up to  $\approx k^2$  triangles, the aggregation around the high  $k$  hubs, visible, e.g., in Fig. 2*a* and *b*, allows the network with a modest number of nodes to absorb a very large number of subgraphs. These calculations can be extended to arbitrary  $(n,m)$  subgraphs, in each case predicting a power law for both  $T(k)$  and  $P(T)$ , with exponents that depend on the  $(n,m)$  parameters (see supporting information). To test the validity of



