

Genomic analysis of regulatory network dynamics reveals large topological changes

Nicholas M Luscombe^{1*+}, M Madan Babu^{2*+}, Haiyuan Yu¹,
Michael Snyder³, Sarah A Teichmann²⁺ and Mark Gerstein^{1,4+}

Department of Molecular Biophysics and Biochemistry¹,
Department of Molecular, Cellular and Developmental Biology³,
Department of Computer Science⁴,
Yale University
PO Box 208114, New Haven, CT 06520-8114, USA

Division of Structural Studies²,
MRC Laboratory of Molecular Biology,
Hills Road, Cambridge CB2 2QH, UK

* These authors contributed equally to this work

+ Correspondence should be addressed to:
Email: sandy@bioinfo.mbb.yale.edu;
Tel: +1 203 432 6105 Fax: +1 360 838 7861

Network analysis has been applied widely, providing a unifying language to describe disparate systems ranging from social interactions to power-grids. It has recently been used in molecular biology, but so far the resulting networks have only been analyzed statically¹⁻⁸. Here we present the dynamics of a biological network on a genomic scale, by integrating transcriptional regulatory information⁹⁻¹¹ and gene-expression data¹²⁻¹⁶ for multiple conditions in *Saccharomyces cerevisiae*. We develop an approach for the *Statistical Analysis of Network Dynamics (SANDY)*, combining well-known global topological measures, local motifs and newly derived statistics. We uncover large changes in underlying network architecture that are unexpected given current viewpoints and random simulations. In response to diverse stimuli, transcription factors (TFs) alter their interactions to varying degrees, thereby rewiring the network. A few TFs serve as permanent hubs, whilst most act transiently only during certain conditions. Looking at sub-network structures, we show environmental responses facilitate fast signal propagation (*eg* with short regulatory cascades), whereas the cell cycle and sporulation direct temporal progression through multiple stages (*eg* with highly inter-connected TFs). Indeed, to drive the latter processes forward, phase-specific TFs inter-regulate serially, and ubiquitously active TFs layer above them in a two-tiered hierarchy. We anticipate many of the concepts presented here – particularly large-scale topological changes and hub transience – will apply to other biological networks, including complex sub-systems in higher eukaryotes.

We begin by assembling a static representation of known regulatory interactions from the results of genetic, biochemical and ChIP-chip experiments. Figure 1 illustrates the complexity of the resultant network, which contains 7,074 regulatory interactions between 142 TFs and 3,420 target genes (interactions can be between TFs and non-TF targets, or two TFs). To get a dynamic perspective, we integrate gene-expression data for five conditions: cell cycle¹³, sporulation¹⁴, diauxic shift¹², DNA damage¹⁶, and stress response¹⁵. From these data, we trace paths in the regulatory network that are active in each condition using a back-tracking algorithm (see Methods).

Figure 1b presents the sub-networks active under different cellular conditions, and gross changes are apparent in the distinct sections of the network that are highlighted. Recent functional genomics studies have analyzed the dynamics of a few TFs^{17,18}; however, Figure 1 represents the first dynamic view of a genome-scale network.

Half of the targets are uniquely expressed in only one condition; in contrast, most TFs are used across multiple processes. The active sub-networks maintain or rewire regulatory interactions, and over half of the active interactions (1,476 of 2,476 total) are completely supplanted by new ones between conditions. Just 66 interactions are retained across four or more conditions; these comprise *hot links*⁶ that are "always on" (compared with the rest of the network) and mostly regulate house-keeping functions.

The large number of changing interactions makes rigorous comparison of active sub-networks impossible visually. Consequently, we introduce an approach called *SANDY* that combines: *standard measures* of network connectivity (involving global topological statistics⁶ and local network motifs⁴), newly derived *follow-on* statistics and comparisons against *simulated controls* to assess the significance of each observation.

Overall, our calculations divide the five condition-specific sub-networks into two categories: *endogenous* and *exogenous* (Figure 1). This allows us to rationalize the different sub-network structures in terms of the biological requirements of each condition. Endogenous processes (cell cycle and sporulation) are multi-stage and operate with an internal transcriptional program, whereas exogenous states (diauxic shift, DNA damage and stress response) constitute binary events that react to external stimuli with a rapid turnover of expressed genes.

We begin *SANDY* by examining global *topological measures* that quantify network architecture (Figure 1c) ⁶. The view from recent studies is that these statistics are remarkably constant across many biological networks (including regulatory systems)^{1,5,6,19,20}. Moreover, most of them remain invariant between randomly simulated sub-graphs of different sizes (Methods).

In fact, we show that topological measures change considerably in the endogenous and exogenous sub-networks. (Furthermore, most of the observed measurements differ significantly from random expectation and are insensitive to addition of noise in the underlying network; Methods). The *in-degree* (k_{in}) is the number of incoming edges per node (*ie* the number of TFs regulating a target). Its average across each sub-network decreases by 20% from endogenous to exogenous conditions. (The probability p that these values originate from the same population is $<3 \times 10^{-4}$; Supplementary Material). The *out-degree* (k_{out}) represents the number of outgoing edges per node (*ie* the number of target genes for each TF). Average values double from endogenous to exogenous conditions ($p < 2 \times 10^{-3}$). The *path length* (l) is the shortest distance between two nodes (*ie* here, it is the number of intermediate regulators between a TF and a terminating target gene). Its average halves from endogenous to exogenous conditions ($p < 10^{-10}$). Finally, the *clustering coefficient* (c) gauges the level of inter-connectivity around a node (*ie* the level of TF inter-regulation). Values range from 0 for totally dispersed nodes to 1 for fully connected ones. Average coefficients nearly halve from endogenous to exogenous conditions ($p < 0.01$).

In biological terms, the small in-degrees for exogenous conditions indicate TFs regulating in simpler combinations, and the large out-degrees signify that each TF has greater regulatory influence by targeting more genes simultaneously. The short paths imply faster propagation of the regulatory signal. Conversely, long paths in the multi-stage, endogenous conditions suggest slower action arising from the formation of regulatory chains to control intermediate phases. Finally, high clustering coefficients in endogenous conditions signify greater inter-regulation between TFs. In summary, sub-networks have evolved to produce rapid, large-scale responses in exogenous states, and carefully coordinated processes in endogenous conditions (Figure 1a).

SANDY also examines sub-networks locally by calculating the occurrence of *network motifs* ⁴, which are compact, specific patterns of inter-connection between TFs and targets. We show the occurrence of the most common in Figure 1c: single-input, multiple-input, and feed-forward-loop motifs (SIMs, MIMs, and FFLs). In SIMs a single TF targets many genes; in MIMs multiple TFs co-regulate sets of genes; and in FFLs a primary TF regulates a secondary one, and both target a final gene. Motifs appear at similar relative frequencies across regulatory networks of diverse organisms (though

individual motifs are not conserved)^{3,7}, and this is also true for the randomly simulated sub-graphs. Therefore, constancy in motif-usage is expected across conditions.

However, Figure 1c shows that the relative occurrence of motifs varies considerably between endogenous and exogenous conditions ($p < 10^{-9}$). SIMs are favoured in exogenous sub-networks where they comprise >55% of regulatory interactions in motifs. But the frequency drops to ~35% in endogenous processes. Instead, these states favour FFLs (~44%). (MIMs do not significantly change their usage).

Previous studies defined precise regulatory properties and information processing tasks for motifs⁴. SIMs and MIMs are implicated in conferring similar regulation over groups of genes, so they are ideal for directing the large-scale gene activation found in exogenous conditions. FFLs are buffers that respond only to persistent input signals. They are suited for endogenous conditions, as cells cannot initiate a new stage until the previous one has stabilized. (FFLs are used sparingly in exogenous processes but may be important in filtering spurious external stimuli).

Having quantified global and local changes with standard topological measures, we now move to the follow-on statistics in *SANDY* (Figure 2). Like many large-scale networks, the regulatory system displays scale-free characteristics (the probability P_k that a TF targets k genes is proportional to $k^{-\gamma}$ for constant γ). This behaviour (maintained across all active sub-networks) signifies the presence of regulatory *hubs* targeting disproportionately large numbers of genes. Hubs are of general interest as they represent the most influential components of a network⁶ and accordingly, tend to be essential²¹. They are considered to target a broad spectrum of gene functions^{4,11,22}, and are commonly located upstream in the network² to expand their influence via secondary TFs²³. These observations suggest that hubs would be invariant features of the network across conditions, and this expectation is supported by the random simulations that converge on similar sets of TF hubs.

Figure 2a shows the observed regulatory hubs in each of the five conditions (Methods). They divide into two groups. The smaller one represents *permanent* hubs, which in line with expectation, are important regardless of cellular state. They mainly comprise multi-functional TFs (eg Abf1) and house-keeping regulators (eg Mig1/2), and are responsible for maintaining hot links. However contrary to expectation, most hubs (78%) are *transient*; they are influential in one condition, but less so in others. Exogenous conditions have fewer hubs, suggesting a more centralized command structure. (This is reflected in different γ , Supplementary Material). About half of the transient hubs are known to be important for their respective conditions (eg Swi4 in the cell cycle; Methods). For the remainder with sparse annotations, their "transient-hub" status in a particular condition considerably augments their functional annotation (eg Sok2 in the cell cycle). Intriguingly, these hubs may also relate to condition-dependent lethality, and this has clear implications for identifying specific drug targets.

The defining feature of transient hubs is their capacity to change interactions between conditions. We attempt to quantify this rewiring more broadly for every TF in the network with the *interchange index*, I . This is defined so that higher values associate with TFs replacing a larger fraction of their interactions. Its histogram reveals a uni-modal central distribution with two groups of extreme outliers (Figure 2b). At one extreme

($I \leq 10\%$), 12 TFs retain all interactions across multiple states. At the other end ($I \geq 90\%$), 27 TFs replace all interactions in switching conditions. Many of these are so extreme that they only regulate genes in a single condition and are inactive otherwise. These include six transient hubs of known importance for the cell cycle and stress response. Most TFs interchange only part of their interactions ($10\% < I < 90\%$). This group comprises most of the hubs; somewhat surprisingly, permanent hubs interchange interactions as often as transient ones, but over a larger number of conditions. Furthermore, TFs in this group often regulate genes of distinct functions in different conditions, so shifting regulatory roles. For example, the permanent hub Abf1 regulates cell growth during endogenous conditions, but refocuses to intracellular transport in stress response (in addition to its maintained core functions).

The rewiring highlighted by the interchange index allows TFs to be active in many conditions. Indeed, 95 of the 142 TFs are used in more than one process (Figure 1b). Specifically, within endogenous conditions 53 of 92 TFs overlap between cell cycle and sporulation (Figure 2c), and there is a similar overlap for exogenous conditions (Supplementary Material). With so much intersection in the repertoire of active TFs, the precise regulation of a condition cannot arise from the specificity of individual TFs. As others have observed²⁴, *combinatorial TF usage* appears to be the key. We calculate that there are 360 unique pair-wise TF combinations (*ie* two TFs regulating the same target) used in at least one condition. In contrast to individual TFs, only a minor proportion of pairs (51 of 360) participate in multiple processes and just 3 of 149 pairs overlap between endogenous conditions (Figure 2b).

Thus far we have focused on the large dynamic changes occurring *between* different cellular conditions. However, dynamic transitions also take place *within* individual processes. Earlier, *SANDY* defined endogenous sub-networks by their long paths and high clustering. We can study the source of these observations by looking at the full scope of inter-regulation between TFs during the cell cycle (Figure 3). This is possible as Cho *et al.*¹³ provide expression-level measurements throughout the cell cycle and assign differentially expressed genes to one of five phases (early G1, late G1, S, G2, M). We then back-tracked from the classified genes to identify active sub-networks during each phase (Methods).

A cluster diagram (Figure 3a) shows that most TFs active in the cell cycle operate only in a particular phase (*eg* Swi4 in late G1). Additionally, a sizeable minority of TFs is ubiquitously active throughout the whole cycle. We uncover two major forms of TF inter-regulation. In *serial inter-regulation*²⁵ (Figure 3b), the phase-specific TFs regulate each other in a sequential manner to drive the cell cycle forward. In fact, we detect complete loops of interactions within the complex circuitry, and the resulting regulatory cascades undoubtedly create the long paths. We also introduce the concept of *parallel inter-regulation* (Figure 3c), where the ubiquitous TFs control the phase-specific ones in a two-tiered system. This effectively provides a stable signal to aid the transition between phases. Furthermore, as about a third of ubiquitous TFs comprise permanent hubs, they may provide a channel of communication to relate the cell-cycle progression with house-keeping functions. Similar observations apply to sporulation (Supplementary Material).

In summary, *SANDY* presents an approach to examine biological network dynamics. In applying it to the yeast regulatory system, it becomes apparent that many observations

made in the static state are not applicable to the condition-specific sub-networks. However in refocusing to a dynamic perspective, we uncover substantial topological changes in network structure, and we capture the essence of the transcriptional regulatory data in a new way. Because of limitations in current datasets, we can examine this only through integrating gene-expression information. However we anticipate future experiments to determine condition-specific interactions directly. Given the robustness of the observations to large perturbations (Methods), we expect our approach and findings to remain valid for these new datasets. Furthermore, we anticipate that many of the concepts we introduce could be readily transferred to other types of biological networks, and complex sub-systems in multi-cellular organisms such as those directing the circadian cycle²⁶ and cellular development.

Methods

Detailed descriptions of the methods are in the Supplementary Material and at <http://sandy.topnet.gersteinlab.org>.

Datasets. (i) The transcriptional regulatory network is assembled from the results of genetic, biochemical and ChIP-chip experiments, with non-DNA-binding factors removed⁹⁻¹¹. (ii) The gene-expression data are compiled from 240 microarray experiments for five conditions¹²⁻¹⁶. We identify the following numbers of genes with differential expression: cell cycle, 455; sporulation, 477; diauxic shift, 1,823; DNA damage, 1,718; and stress response, 866.

Back-tracking algorithm. This defines sections of the regulatory network used in each condition: (i) We identify TFs *present* in a condition as those with sufficiently high expression levels. (ii) We flag differentially expressed genes that appear in the regulatory network. (iii) We mark as *active* the regulatory links between present TFs and differentially expressed genes. (iv) We then search for any other present TFs that are linked to a TF with an already active link and make this connection active. The last step is repeated until no more links are made active. The same procedure identifies sub-networks active in particular phases of cell cycle and sporulation.

SANDY. This extends the methodology used by the TopNet software tool²⁷ and it evaluates each sub-network with the following. (i) Standard statistics including global measures of topology (k_{in} , k_{out} , l , c)⁶ and local motif occurrence (SIM, MIM, FFL)⁴. (ii) Follow-on statistics including: (a) permanent and transient hub identification, (b) interchange index (I), (c) counting the overlap in TF usage (individual and pairs) across multiple conditions. (Hubs are TFs in the top 30%, by number of target genes, in at least one condition. The number of target genes is normalized to measure the relative influence of a TF hub in a particular process). In all cases, regulatory functions are obtained from SGD²⁸ and are current as of June 2004. (iii) We compare observations with random expectation by simulating sub-graphs that are similar in size to each sub-network, and calculating statistics (i) and (ii) for them. Simulated sub-graphs sample the same number of “differentially expressed” genes and back-track through the static network. We also test the sensitivity of our observations to noise by randomly perturbing the static networks by 30% (random addition, deletion and replacement of interactions), back-tracking from the original differentially expressed genes, and then recalculating the statistics.

Declaration of competing financial interests

The authors declare that they have no competing financial interests.

Acknowledgements

We thank Paul Bertone, Nuria Domedel Puig, Eivind Hovig, Ronald Jansen, Kristine Kleivi, Georgy Koentges, Eugene Koonin, Boris Lenhard, Alberto Paccanaro, Joel Rozowsky, Jesper Tegner and Annabel Todd for insightful comments on the paper. NML thanks the Anna Fuller Fund and the MRC LMB Visitor's Program. MMB acknowledges financial support from the Cambridge Commonwealth Trust, Trinity College, Cambridge and the MRC LMB. MG is supported by the NIH.

References

1. Jeong, H., Tombor, B., Albert, R., Oltvai, Z. N. & Barabasi, A. L. The large-scale organization of metabolic networks. *Nature* **407**, 651-4 (2000).
2. Guelzim, N., Bottani, S., Bourguine, P. & Kepes, F. Topological and causal structure of the yeast transcriptional regulatory network. *Nat Genet* **31**, 60-3 (2002).
3. Milo, R. et al. Network motifs: simple building blocks of complex networks. *Science* **298**, 824-7 (2002).
4. Shen-Orr, S. S., Milo, R., Mangan, S. & Alon, U. Network motifs in the transcriptional regulation network of Escherichia coli. *Nat Genet* **31**, 64-8 (2002).
5. Oltvai, Z. N. & Barabasi, A. L. Systems biology. Life's complexity pyramid. *Science* **298**, 763-4 (2002).
6. Barabasi, A. L. & Oltvai, Z. N. Network biology: understanding the cell's functional organization. *Nat Rev Genet* **5**, 101-13 (2004).
7. Milo, R. et al. Superfamilies of evolved and designed networks. *Science* **303**, 1538-42 (2004).
8. Teichmann, S. A. & Babu, M. M. Gene regulatory network growth by duplication. *Nat Genet* **36**, 492-6 (2004).
9. Svetlov, V. V. & Cooper, T. G. Review: compilation and characteristics of dedicated transcription factors in *Saccharomyces cerevisiae*. *Yeast* **11**, 1439-84 (1995).
10. Horak, C. E. et al. Complex transcriptional circuitry at the G1/S transition in *Saccharomyces cerevisiae*. *Genes Dev* **16**, 3017-33 (2002).
11. Lee, T. I. et al. Transcriptional regulatory networks in *Saccharomyces cerevisiae*. *Science* **298**, 799-804 (2002).
12. DeRisi, J. L., Iyer, V. R. & Brown, P. O. Exploring the metabolic and genetic control of gene expression on a genomic scale. *Science* **278**, 680-6 (1997).
13. Cho, R. J. et al. A genome-wide transcriptional analysis of the mitotic cell cycle. *Mol Cell* **2**, 65-73 (1998).
14. Chu, S. et al. The transcriptional program of sporulation in budding yeast. *Science* **282**, 699-705 (1998).
15. Gasch, A. P. et al. Genomic expression programs in the response of yeast cells to environmental changes. *Mol Biol Cell* **11**, 4241-57 (2000).

16. Gasch, A. P. et al. Genomic expression responses to DNA-damaging agents and the regulatory role of the yeast ATR homolog Mec1p. *Mol Biol Cell* **12**, 2987-3003 (2001).
17. Odom, D. T. et al. Control of pancreas and liver gene expression by HNF transcription factors. *Science* **303**, 1378-81 (2004).
18. Zeitlinger, J. et al. Program-specific distribution of a transcription factor dependent on partner transcription factor and MAPK signaling. *Cell* **113**, 395-404 (2003).
19. Watts, D. J. & Strogatz, S. H. Collective dynamics of 'small-world' networks. *Nature* **393**, 440-2 (1998).
20. Wagner, A. & Fell, D. A. The small world inside large metabolic networks. *Proc R Soc Lond B Biol Sci* **268**, 1803-10 (2001).
21. Yu, H., Greenbaum, D., Xin Lu, H., Zhu, X. & Gerstein, M. Genomic analysis of essentiality within protein networks. *Trends Genet* **20**, 227-31 (2004).
22. Martinez-Antonio, A. & Collado-Vides, J. Identifying global regulators in transcriptional regulatory networks in bacteria. *Curr Opin Microbiol* **6**, 482-9 (2003).
23. Madan Babu, M. & Teichmann, S. A. Evolution of transcription factors and the gene regulatory network in Escherichia coli. *Nucleic Acids Res* **31**, 1234-44 (2003).
24. Pilpel, Y., Sudarsanam, P. & Church, G. M. Identifying regulatory networks by combinatorial analysis of promoter elements. *Nat Genet* **29**, 153-9 (2001).
25. Simon, I. et al. Serial regulation of transcriptional regulators in the yeast cell cycle. *Cell* **106**, 697-708 (2001).
26. Ueda, H. R. et al. A transcription factor response element for gene expression during circadian night. *Nature* **418**, 534-9 (2002).
27. Yu, H., Zhu, X., Greenbaum, D., Karro, J. & Gerstein, M. TopNet: a tool for comparing biological sub-networks, correlating protein properties with topological statistics. *Nucleic Acids Res* **32**, 328-37 (2004).
28. Christie, K. R. et al. Saccharomyces Genome Database (SGD) provides tools to identify and analyze sequences from Saccharomyces cerevisiae and related sequences from other organisms. *Nucleic Acids Res* **32**, D311-4 (2004).

Figures

Figure 1. Dynamic representation of the transcriptional regulatory network and standard statistics. (a) Schematics and summary of properties for the endogenous and exogenous sub-networks. (b) Graphs of the static and condition-specific networks. TFs and target genes are shown as nodes in the upper and lower sections of each graph respectively, and regulatory interactions are drawn as edges; they are coloured by the number of conditions in which they are active. Different conditions use distinct sections of the network. (c) Standard statistics (global topological measures and local network motifs) describing network structures. These vary between endogenous and exogenous conditions; those that are high compared with other conditions are shaded. (Note, the graph for the static state displays only sections that are active in at least one condition, but the table provides statistics for the entire network including inactive regions).

Figure 2. Derived "follow-on" statistics for network structures. (a) TF hub usage in different cellular conditions. The cluster diagram shades cells by the normalized number of genes targeted by TF hubs in each condition. One cluster represents permanent hubs and the others condition-specific transient hubs. Genes are labelled with four-letter names when they have an obvious functional role in the condition, and seven-letter ORF names when there is no obvious role. Of the latter, gene names are red and italicised when annotation is sparse. Starred hubs show extreme interchange index values, $I = 1$. (b) Interaction interchange (I) of TFs between conditions. A histogram of I for all active TFs shows a uni-modal distribution with two extremes. Pie charts show five example TFs with different proportions of interchanged interactions. We list the main functions of the distinct target genes regulated by each example TF. Note how the TFs' regulatory functions change between conditions. (c) Overlap in TF usage between conditions. Venn diagrams show the numbers of individual TFs (large intersection) and pair-wise TF combinations (small intersection) that overlap between the two endogenous conditions.

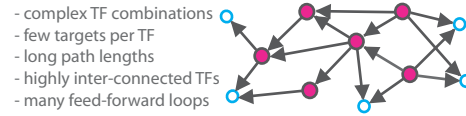
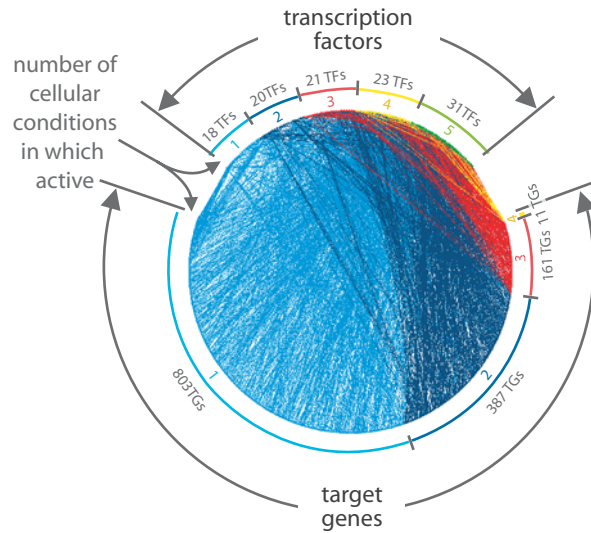
Figure 3. TF inter-regulation during the cell cycle. (a) The 70 TFs active in the cell cycle. The diagram shades each cell by the normalized number of genes targeted by each TF in a phase. Five clusters represent phase-specific TFs and one cluster is for ubiquitously active TFs. TF names are given in Supplementary Material. (b) Serial inter-regulation between phase-specific TFs. Network diagrams show that TFs active in one phase regulate TFs in subsequent phases. In the late phases, TFs appear to regulate those in the next cycle. (c) Parallel inter-regulation between phase-specific and ubiquitous TFs in a two-tiered hierarchy. Serial and parallel inter-regulation operate in tandem to drive the cell cycle while balancing it with basic house-keeping processes.

static

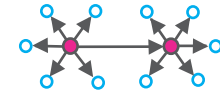
endogenous

exogenous

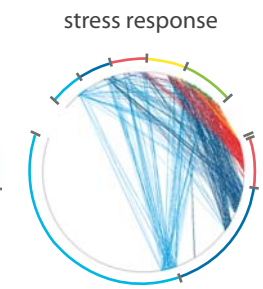
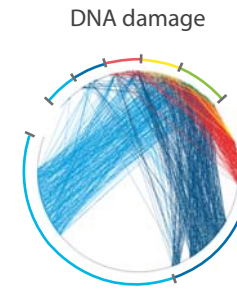
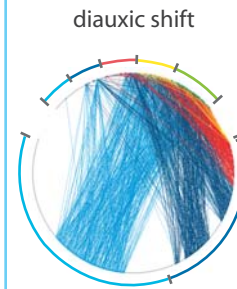
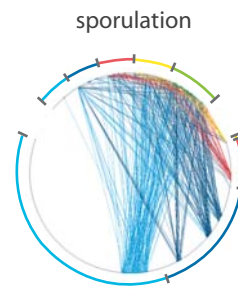
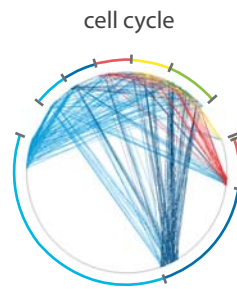
(a) schematic and summary



- simple TF combinations
- many targets per TF
- short path lengths
- few inter-connected TFs
- many single input motifs



(b) graph



(c) standard network statistics

size	static	cell cycle	sporulation	diauxic shift	DNA damage	stress response
# transcription factors	142	70	74	71	72	63
# target genes	3,420	280	257	748	678	362
# regulatory interactions	7,074	550	481	1,217	1,082	566
in-degree ($\langle k_{in} \rangle$)	2.1	2.0	1.9	1.6	1.6	1.6
out-degree ($\langle k_{out} \rangle$)	49.8	7.9	6.5	17.1	15.0	9.0
path length ($\langle l \rangle$)	4.7	4.5	3.4	2.1	2.0	2.2
clustering coefficient ($\langle c \rangle$)	0.11	0.15	0.14	0.09	0.09	0.08
single input (SIM)	1,748 (37.6%)	130 (32.0%)	117 (38.9%)	438 (57.4%)	462 (55.7%)	228 (59.1%)
multiple input (MIM)	325 (7.0%)	96 (23.7%)	50 (16.6%)	180 (23.6%)	226 (27.3%)	78 (20.2%)
feed-forward loop (FFL)	2,581 (55.5%)	180 (44.3%)	134 (44.5%)	145 (19.0%)	141 (17.0%)	80 (20.7%)
total	4,654	406	301	763	829	386

