

Genomes, phylogeny, and evolutionary systems biology

Mónica Medina*

Department of Evolutionary Genomics, Department of Energy Joint Genome Institute, 2800 Mitchell Drive, Walnut Creek, CA 94598

With the completion of the human genome and the growing number of diverse genomes being sequenced, a new age of evolutionary research is currently taking shape. The myriad of technological breakthroughs in biology that are leading to the unification of broad scientific fields such as molecular biology, biochemistry, physics, mathematics, and computer science are now known as systems biology. Here, I present an overview, with an emphasis on eukaryotes, of how the postgenomics era is adopting comparative approaches that go beyond comparisons among model organisms to shape the nascent field of evolutionary systems biology.

postgenomics | eukaryotic genomics | network evolution

Systems biology is in the eye of the beholder.
Leroy Hood

Only in the last decade have we had access to nearly complete genomes of a diversity of organisms allowing for large-scale comparative analysis. The access to this immense amount of data is providing profound insight into the tree of life at all levels of divergence (Fig. 1*A*). It is thus not surprising that understanding phylogenetic relationships is a prevalent research goal among not only evolutionary biologists but also all scientists interested in the organization and function of the genome. New genome sequences and analysis methods are helping improve our understanding of phylogeny, and at the same time improved phylogenies and phylogenetic theory are generating a better understanding of genome evolution. Currently however, the level of genome sequencing for different branches of the tree of life is far from equivalent. Prokaryotic genome projects are abundant, mainly due to their small genome sizes, with >200 genomes already published and at least 500 currently in progress (www.genomesonline.org). In contrast, <300 eukaryotic genomes are either finished or in progress (www.genomesonline.org). Nevertheless, these data are starting to have a major impact on our understanding of eukaryotic evolution.

These new genomic data have informed our understanding of phylogenetic relationships, and the emerging consensus topologies are adding new insight to the small subunit ribosomal RNA phylogenies. For example, the topology of the ribosomal eukaryotic tree has been recently redrawn with the use of genomic signatures that place the root of all eukaryotic life between two newly uncovered major clades, Unikonts and Bikonts (Fig. 1*A*). Unikonts, which contain the heterotrophic groups Opisthokonta and the Amoebozoa, share a derived three-gene fusion of enzyme-encoding genes in the pyrimidine synthesis pathway (1), whereas Bikonts, which contain the remaining eukaryotic clades, share another derived gene fusion between dihydrofolate reductase and thymidine synthase (2). All photosynthetic groups of primary and secondary plastid symbiotic origins are now thought to be within the Bikonts. Although the animal, fungal, and plant lineages are the most widely represented in terms of genome initiatives (Fig. 1*B–D*), it is significant that multiple protistan

genome projects have also been initiated by the interest of diverse scientific communities, including parasitologists (3), plant pathologists (4), oceanographers (5), and evolutionary biologists (www.biology.uiowa.edu/workshop).

As more whole-genome projects are being completed, post-genomic biology is also providing insight into the function of biological systems by the use of new high-throughput bioanalytical methods, information technology, and computational modeling. This new revolution in biology has become known as systems biology (6). In addition to shifting approaches to biological research from reductionist strategies to pathway- and system-level strategies (7), another paradigm is rapidly emerging, namely the use of phylogenetically based inference in systems biology. Before the genomic revolution, research questions were typically addressed within a single model organism, with only occasional comparative studies when similar information was available for another organism. These comparisons were made between distantly related taxa, and the evolutionary implications were rarely mentioned or taken into account. The increasing importance of comparative analysis is evident in the growing proportion of new prokaryotic genome projects that have been chosen primarily because of their phylogenetic relationship to model organisms, such as *Escherichia coli* and *Bacillus subtilis* and their corresponding related taxa. This same trend is occurring for eukaryotes. Some prominent examples are the multiple *Saccharomyces* genome projects and those of other ascomycete fungi, the several *Plasmodium* projects and other genome initiatives for apicomplexan taxa, the numerous *Caenorhabditis* and other nematode genome projects, the multiple *Drosophila* and arthropod genome projects, and the large number of primate and mammalian genome projects.

Genomes and Phylogeny of Higher Eukaryotes

Metazoa. The sampling of the metazoan tree, and in particular of the chordate branch, was undertaken primarily due to the usefulness of the genomes in understanding human biology. However, this larger genomic dataset is already providing a powerful tool for comparative analysis and more accurate evolutionary inference. Deeper divergences in the Metazoan tree have become the target of major scrutiny due to the interest in comparative developmental genetics (Fig. 1*B*). Based on molecular phylogenies, the bilaterian phyla have been rearranged into three large clades, deuterostomes, lophotrochozoans, and ecdysozoans, these last two being sister taxa inside the protostome clade. At present, there is still debate regarding the placement of nematodes in the tree (i.e., the Ecdysozoa vs. Coelomata hypotheses) because analysis of genomic data

This paper results from the Arthur M. Sackler Colloquium of the National Academy of Sciences, "Systematics and the Origin of Species: On Ernst Mayr's 100th Anniversary," held December 16–18, 2004, at the Arnold and Mabel Beckman Center of the National Academies of Science and Engineering in Irvine, CA.

*Present address: School of Natural Sciences, University of California, P.O. Box 2039, Merced, CA 95344. E-mail: mmedina@ucmerced.edu.

© 2005 by The National Academy of Sciences of the USA

networks such as their “scale-free” nature and “small world” organization. In scale-free networks, a few nodes (hubs) have the largest number of connections to other nodes, whereas most of the nodes have just a few connections. This property is reflected in a power-law distribution. In practical terms, this relationship means that, in a protein interaction network, most proteins interact with a couple of others whereas a few proteins (hubs) interact with a large number, and that, in a metabolic network, a few molecules (hubs) participate in most reactions whereas the rest participate in one or two. The “small world” concept refers to the property of such spoke-and-hub networks that there is a small path length between nodes, just as in modern air travel where only a few flights connect any two cities in the world. This property means that a path of just a few interactions or reactions will connect almost any pair of molecules in the cell (29).

Additional levels in the hierarchy of biological networks and the interactions between them are now being characterized that will allow for integration of data and new theoretical predictions (32). Processes widely studied by evolutionary biologists such as selection, gene duplication, and neutral evolution are being examined in the context of network models as opposed to at the level of individual genes or molecules (33–37).

Evolution of Biological Networks

Transcriptional Networks. High-throughput global gene expression approaches such as EST sequencing and microarrays are now common practice for functional assessment of the genome. The extensive microarray gene expression datasets available for model and non-model organisms are starting to be incorporated into a comparative approach to study transcriptome evolution at multiple levels of divergence. At lower levels of divergence, studies in organisms including fish (38), fruitfly (39, 40), and yeast (41) have now shown that extensive variation exists in the transcriptome in natural populations and that this variation is likely to be an important factor in organismal evolution. Transcriptome comparisons across several primate and mouse species, however, suggest that the majority of gene expression differences within and between species evolve in a selectively neutral or nearly neutral fashion (42). At intermediate levels of divergence, less information is available at present due to lack of genomic data. Although analytically challenging, the use of gene expression profiling by heterologous hybridization to a single species cDNA microarray is starting to be explored, potentially opening the door to comparative analyses of taxa as divergent as 200 mega-annum (Ma) (43). This application would be of great significance for the comparative study of non-model organisms that are only distantly related to an already sequenced species. At deep levels of divergence, coexpression of large aggregates of functionally related genes seems to be conserved across evolution. Two recent comparisons of the transcriptomes of several of the model organisms [*S. cerevisiae*, *D. melanogaster*, *C. elegans*, and *H. sapiens* in one case (44), and these four plus *A. thaliana* and *E. coli* in the second case (45)] support the hypothesis that coexpression networks can be split into multiple components enriched for genes involved in similar functional processes. Some of these identified components can be unique to a certain clade, such as the signaling pathway and neuronal function components present only in metazoans in the four-species comparison (44). These cross-species comparisons promise to provide more information about coexpression network evolution as the transcriptomes of additional diverse lineages becomes available (46).

Central to postgenomic analysis is the accuracy of genome annotation. The degree of accuracy in which genomes are annotated is affected by the quality of sequence assembly, gene prediction, and functional annotation by both bioinformatics and experimental data. This relationship is particularly critical in genome projects of non-model organisms where little genetic work has been performed in the past. All these factors, combined

with the lack of network information outside the model organisms, point to the tradeoff between a comprehensive systems analysis of a particular network within a well-studied organism, versus the historical perspective introduced by evolutionary conservation or divergence of systems through time in phylogenetic comparisons. Therefore, although only partial inference is possible at present, studies have already shown that the comparative approach to coexpression not only is giving insight into the universal rules that govern biological systems but also has practical implications by helping improve functional annotations of both model and non-model organisms (44, 45). Because comparative analyses of coexpression data from several model organisms have shown high levels of conservation between such divergent taxa as prokaryotes (*E. coli* and *B. subtilis*) (47), opisthokont eukaryotes (44), and even prokaryotes and eukaryotes (45), some efforts are now targeting the coupled evolution of regulatory networks and the transcriptome.

Regulatory Networks. The characterization of the transcriptome is only a fraction of the information needed to understand global cellular processes because gene expression is driven by the spatio-temporal localization of regulatory networks and details of specific protein–DNA and protein–protein interactions. Genome-wide efforts to characterize transcriptional regulatory networks have already been fruitful in model organisms like yeast (48) and *E. coli* (49). In multicellular organisms, fractions of the regulatory networks are being characterized for sea urchins (50), *Drosophila* (51), and mammals (52).

Transcription factors are regulatory proteins that influence the expression of specific genes. They work by binding to cis-regulatory elements (short and often degenerate sequence motifs frequently located upstream of genes) where they interact with the transcription apparatus to either enhance or repress gene expression. Even though identifying cis-regulatory elements in new genomes is an inherently difficult task due to their short sequence length and as yet unknown syntax, comparative approaches have been helpful. By aligning orthologous regions flanking a gene from multiple species, conserved noncoding sequence motifs can be distinguished. These evolutionary conserved motifs are then hypothesized to be potential functional elements. This method, called phylogenetic footprinting (53), has successfully been used to identify a limited number of regulatory regions in vertebrates (54, 55) and plants (56, 57). More sophisticated comparative approaches are starting to combine computational prediction and laboratory validation of regulatory networks. Coexpression data and known cis-regulatory elements from *S. cerevisiae* were used in a multispecies comparison of 13 published ascomycete genomes, finding multiple cases of regulatory conservation but also some cases of regulatory diversification (58). It has become apparent, however, that sequence conservation alone will not help identify all cis-regulatory elements by phylogenetic footprinting, and additional data and experimental approaches have to be integrated (59).

Gene expression can be regulated not only at transcriptional initiation but also at other levels, such as during mRNA editing, transport, or translation, and characterizing these interactions and their evolution is one of the many future challenges of systems biology (60). For example, comparative work on populations of yeast and fruitfly has recently shown that protein–protein interactions are negatively associated with evolutionary variation in gene expression (61). A comparative analysis of the *E. coli* and yeast regulatory networks has demonstrated that gene duplication has a key role in network evolution both in eukaryotes and prokaryotes (62). Finally, introducing concepts of network dynamics has revealed new topological changes in the regulatory network in yeast (63), an approach that, incorporated into a comparative framework, will eventually provide answers

19. Martinez, D., Larrondo, L. F., Putnam, N., Gelpke, M. D., Huang, K., Chapman, J., Helfenbein, K. G., Ramaiya, P., Detter, J. C., Larimer, F., *et al.* (2004) *Nat. Biotechnol.* **22**, 695–700.
20. Pryer, K. M., Schneider, H., Zimmer, E. A. & Ann Banks, J. (2002) *Trends Plant Sci.* **7**, 550–554.
21. Archibald, J. M. & Keeling, P. J. (2002) *Trends Genet.* **18**, 577–584.
22. Arabidopsis Genome Initiative (2000) *Nature* **408**, 796–815.
23. Yu, J., Hu, S., Wang, J., Wong, G. K., Li, S., Liu, B., Deng, Y., Dai, L., Zhou, Y., Zhang, X., *et al.* (2002) *Science* **296**, 79–92.
24. Goff, S. A., Ricke, D., Lan, T. H., Presting, G., Wang, R., Dunn, M., Glazebrook, J., Sessions, A., Oeller, P., Varma, H., *et al.* (2002) *Science* **296**, 92–100.
25. Fiehn, O. (2002) *Plant Mol. Biol.* **48**, 155–171.
26. Iderker, T., Galitski, T. & Hood, L. (2001) *Annu. Rev. Genomics Hum. Genet.* **2**, 343–372.
27. Kitano, H. (2002) *Science* **295**, 1662–1664.
28. Wolfe, K. H. & Li, W. H. (2003) *Nat. Genet.* **33**, Suppl., 255–265.
29. Barabasi, A. L. & Oltvai, Z. N. (2004) *Nat. Rev. Genet.* **5**, 101–113.
30. Jeong, H., Tombor, B., Albert, R., Oltvai, Z. N. & Barabasi, A. L. (2000) *Nature* **407**, 651–654.
31. Ravasz, E., Somera, A. L., Mongru, D. A., Oltvai, Z. N. & Barabasi, A. L. (2002) *Science* **297**, 1551–1555.
32. Ge, H., Walhout, A. J. & Vidal, M. (2003) *Trends Genet.* **19**, 551–560.
33. Wagner, A. (September 30, 2003) *Sci. STKE*, 10.1126/stke.2003.202.pe41.
34. Wagner, A. (2003) *Proc. R. Soc. London Ser. B Biol. Sci.* **270**, 457–466.
35. van Noort, V., Snel, B. & Huynen, M. A. (2004) *EMBO Rep.* **5**, 280–284.
36. Wuchty, S. (2004) *Genome Res.* **14**, 1310–1314.
37. Hahn, M. W., Conant, G. C. & Wagner, A. (2004) *J. Mol. Evol.* **58**, 203–211.
38. Olesiak, M. J., Churcill, G. A. & Crawford, D. L. (2002) *Nat. Genet.* **32**, 261–266.
39. Meiklejohn, C. D., Parsch, J., Ranz, J. M. & Hartl, D. L. (2003) *Proc. Natl. Acad. Sci. USA* **100**, 9894–9899.
40. Ranz, J. M., Castillo-Davis, C. I., Meiklejohn, C. D. & Hartl, D. L. (2003) *Science* **300**, 1742–1745.
41. Townsend, J. P., Cavalieri, D. & Hartl, D. L. (2003) *Mol. Biol. Evol.* **20**, 955–963.
42. Khaitovich, P., Weiss, G., Lachmann, M., Hellmann, I., Enard, W., Muetzel, B., Wirkner, U., Ansong, W. & Paabo, S. (2004) *PLoS Biol.* **2**, E132.
43. Renn, S. C., Aubin-Horth, N. & Hofmann, H. A. (2004) *BMC Genomics* **5**, 42.
44. Stuart, J. M., Segal, E., Koller, D. & Kim, S. K. (2003) *Science* **302**, 249–255.
45. Bergmann, S., Ihmels, J. & Barkai, N. (2004) *PLoS Biol.* **2**, E9.
46. Zhou, X. J. & Gibson, G. (2004) *Genome Biol.* **5**, 232.
47. Snel, B., van Noort, V. & Huynen, M. A. (2004) *Nucleic Acids Res.* **32**, 4725–4731.
48. Lee, T. I., Rinaldi, N. J., Robert, F., Odom, D. T., Bar-Joseph, Z., Gerber, G. K., Hannett, N. M., Harbison, C. T., Thompson, C. M., Simon, I., *et al.* (2002) *Science* **298**, 799–804.
49. Shen-Orr, S. S., Milo, R., Mangan, S. & Alon, U. (2002) *Nat. Genet.* **31**, 64–68.
50. Davidson, E. H. (2001) *Genomic Regulatory Systems: Development and Evolution* (Academic, San Diego).
51. Berman, B. P., Pfeiffer, B. D., Laverty, T. R., Salzberg, S. L., Rubin, G. M., Eisen, M. B. & Celniker, S. E. (2004) *Genome Biol.* **5**, R61.
52. ENCODE Project Consortium (2004) *Science* **306**, 636–640.
53. Tagle, D. A., Koop, B. F., Goodman, M., Slightom, J. L., Hess, D. L. & Jones, R. T. (1988) *J. Mol. Biol.* **203**, 439–455.
54. Gumucio, D. L., Heilstedt-Williamson, H., Gray, T. A., Tarle, S. A., Shelton, D. A., Tagle, D. A., Slightom, J. L., Goodman, M. & Collins, F. S. (1992) *Mol. Cell. Biol.* **12**, 4919–4929.
55. Dermitzakis, E. T., Reymond, A., Scamuffa, N., Ucla, C., Kirkness, E., Rossier, C. & Antonarakis, S. E. (2003) *Science* **302**, 1033–1035.
56. Hong, R. L., Hamaguchi, L., Busch, M. A. & Weigel, D. (2003) *Plant Cell* **15**, 1296–1309.
57. Kaplinsky, N. J., Braun, D. M., Penterman, J., Goff, S. A. & Freeling, M. (2002) *Proc. Natl. Acad. Sci. USA* **99**, 6147–6151.
58. Gasch, A. P., Moses, A. M., Chiang, D. Y., Fraser, H. B., Berardini, M. & Eisen, M. B. (2004) *PLoS Biol.* **2**, e398.
59. Richards, S., Liu, Y., Bettencourt, B. R., Hradecky, P., Letovsky, S., Nielsen, R., Thornton, K., Hubisz, M. J., Chen, R., Meisel, R. P., *et al.* (2005) *Genome Res.* **15**, 1–18.
60. Wei, G. H., Liu, D. P. & Liang, C. C. (2004) *Biochem. J.* **381**, 1–12.
61. Lemos, B., Meiklejohn, C. D. & Hartl, D. L. (2004) *Nat. Genet.* **36**, 1059–1060.
62. Teichmann, S. A. & Babu, M. M. (2004) *Nat. Genet.* **36**, 492–496.
63. Luscombe, N. M., Babu, M. M., Yu, H., Snyder, M., Teichmann, S. A. & Gerstein, M. (2004) *Nature* **431**, 308–312.
64. Howard, M. L. & Davidson, E. H. (2004) *Dev. Biol.* **271**, 109–118.
65. Giot, L., Bader, J. S., Brouwer, C., Chaudhuri, A., Kuang, B., Li, Y., Hao, Y. L., Ooi, C. E., Godwin, B., Vitols, E., *et al.* (2003) *Science* **302**, 1727–1736.
66. Li, S., Armstrong, C. M., Bertin, N., Ge, H., Milstein, S., Boxem, M., Vidalain, P. O., Han, J. D., Chesneau, A., Hao, T., *et al.* (2004) *Science* **303**, 540–543.
67. Uetz, P., Giot, L., Cagney, G., Mansfield, T. A., Judson, R. S., Knight, J. R., Lockshon, D., Narayan, V., Srinivasan, M., Pochart, P., *et al.* (2000) *Nature* **403**, 623–627.
68. Lehner, B. & Fraser, A. G. (2004) *Genome Biol.* **5**, R63.
69. Yu, H., Luscombe, N. M., Lu, H. X., Zhu, X., Xia, Y., Han, J. D., Bertin, N., Chung, S., Vidal, M. & Gerstein, M. (2004) *Genome Res.* **14**, 1107–1118.
70. Stitt, M. & Fernie, A. R. (2003) *Curr. Opin. Biotechnol.* **14**, 136–144.
71. Oksman-Caldentey, K. M., Inze, D. & Oresic, M. (2004) *Proc. Natl. Acad. Sci. USA* **101**, 9949–9950.
72. Castrillo, J. I. & Oliver, S. G. (2004) *J. Biochem. Mol. Biol.* **37**, 93–106.
73. Papin, J. A., Reed, J. L. & Palsson, B. O. (2004) *Trends Biochem. Sci.* **29**, 641–647.
74. Walhout, A. J., Reboul, J., Shtanko, O., Bertin, N., Vaglio, P., Ge, H., Lee, H., Doucette-Stamm, L., Gunsalus, K. C., Schetter, A. J., *et al.* (2002) *Curr. Biol.* **12**, 1952–1958.
75. Ge, H., Liu, Z., Church, G. M. & Vidal, M. (2001) *Nat. Genet.* **29**, 482–486.
76. Levesque, M. P. & Benfey, P. N. (2004) *Curr. Biol.* **14**, R179–80.
77. Benfey, P. N. (2004) *Dev. Cell* **7**, 329–330.
78. Thomas, M. A. & Klaper, R. (2004) *Trends Ecol. Evol.* **19**, 439–445.
79. Wake, M. L. (2003) *Integr. Comp. Biol.* **43**, 239–241.
80. Simpson, A. G. & Roger, A. J. (2004) *Curr. Biol.* **14**, R693–R696.
81. Bhattacharya, D., Yoon, H. S. & Hackett, J. D. (2004) *BioEssays* **26**, 50–60.
82. Adoutte, A., Balavoine, G., Lartillot, N. & de Rosa, R. (1999) *Trends Genet.* **15**, 104–108.
83. Medina, M., Collins, A. G., Silberman, J. D. & Sogin, M. L. (2001) *Proc. Natl. Acad. Sci. USA* **98**, 9707–9712.
84. Ruiz-Trillo, I., Paps, J., Loukota, M., Ribera, C., Jondelius, U., Bagaña, J. & Ruitort, M. (2002) *Proc. Natl. Acad. Sci. USA* **99**, 11246–11251.