

Practical computational approaches to inferring protein function

Edward M. Marcotte

A quick search through high-throughput proteomics and genomics data can reveal information on many aspects of protein function, such as mutant phenotypes, protein interactions, mRNA expression patterns, transcriptional regulation, and even protein structure. The computational integration of such data is proving to be the most effective route to protein function.

Department of Chemistry and Biochemistry, Institute for Cellular and Molecular Biology, Center for Computational Biology and Bioinformatics,
2500 Speedway, MBB 3.232,
University of Texas at Austin,
Austin, TX 78712, USA
e-mail: marcotte@icmb.utexas.edu

As the availability of raw data about protein function grows continuously, investigators are scrambling to convert these data to knowledge. Datasets describing deletion mutant phenotypes, protein and mRNA expression profiles, genome sequences, and protein interactions, to name but a few, have opened several new routes to protein function. Each set of data, typically collected on as large a scale as is practical, tells something of the functions of many proteins. As a consequence, few of the current computational approaches for inferring protein function derive from first-principles models of protein function; rather, they represent varied approaches for ‘mining’ functional inferences from diverse proteomics and functional-genomics data. The best routes to function involve integrating the partial functional inferences from many types of these data at once.

[Approaches for the integration of different types of data relating to protein function](#) generally take what might be termed the ‘genome-down’ rather than ‘protein-up’ approach. They exploit the principle that a protein’s function can be determined more easily in the context of the other proteins with which it works in the cell. Genome-down approaches systematically analyze the entire set of proteins encoded by a genome, and only focus in on specific proteins after

completing this holistic analysis. One of the most compelling arguments in support of this strategy is that the proteins with known function act as cases to test the approach’s effectiveness, and the overall accuracy of the approach can be measured. [Such assessments of accuracy show](#) that these methods are often still a bit hit-and-miss, and there is no guarantee that data will exist for a particular protein. Despite these caveats, a tremendous amount of functional information has been found in this manner, much freely available for public consumption. This review discusses recent computational approaches for inferring protein function, several successful integrated approaches for analyzing both functional genomics and proteomics data, and tools for effectively navigating these complex datasets. A flowchart describing the use of these tools is provided in Fig. 1; internet links to the major resources available at present are listed in Boxes 1–3.

Protein function from comparative genomics

Two major trends are emerging in the use of genomics data to infer protein function. The first approach relies upon discovering the information about gene function that is intrinsic in genomes. This information can be revealed by finding contextual cues shared by genes that interact or perform a given function [1,2]. The second approach, which relies upon the completeness of genome sequences, is to match a gene with its equivalent genes in well-characterized model organisms. This approach allows the investigator to profit from the rich functional datasets that exist for model organisms.

When using the first approach, several contextual trends have proved to be useful for finding protein function. These include searching for evidence of fusions between the gene of interest and other functionally related genes [3–5]; finding functionally linked genes because of their tendency to be ‘co-inherited’ [6–9]; identifying proteins that physically interact [by looking for](#) the conservation of their phylogenetic tree structures [10–12]; and

Box 1. Servers for computationally predicting protein function

Prediction of protein function, interactions, and networks

Bioverse	http://bioverse.compbio.washington.edu
<i>In silico</i> two hybrid	http://www.pdg.cnb.uam.es/i2h
InterDom	http://InterDom.lit.org.sg
Magic	http://genome-www.stanford.edu/magic
Predictome	http://predictome.bu.edu
ProtFun	http://www.cbs.dtu.dk/services/ProtFun
ProteinFunction	http://www.aber.ac.uk/compsci/Research/bio/ProteinFunction
Protein Link Explorer (PLEX)	http://bioinformatics.icmb.utexas.edu/plex
STRING	http://www.bork.embl-heidelberg.de/STRING

Predicting prokaryotic operons to find functionally linked proteins

Gene Neighbors	http://bioinformatics.icmb.utexas.edu/operons
STRING	http://www.bork.embl-heidelberg.de/STRING
TUpredictions	http://www.cifn.unam.mx/moreno/pub/TUpredictions
WIT	http://wit.mcs.anl.gov/WIT2

107 computationally identifying operons, either by virtue
108 of their conservation across organisms [13–15] or the
109 physical separation between the genes along the
110 chromosome [16]. Each of these approaches produces
111 a set of candidate proteins that are functionally
112 linked to a protein of interest, with a score that
113 indicates the confidence of the linkages.

114 Comparisons of phylogenetic trees or analyses of
115 the numbers of nucleotides that separate genes can
116 be performed easily. The practical implementation of
117 the remaining methods, however, requires the
118 systematic comparison of large sets of protein
119 sequences from many organisms, followed by
120 statistical analysis of the many comparisons. For
121 example, calculation of phylogenetic-profile-based
122 linkages involves comparing the amino-acid
123 sequence of each protein encoded by a genome with
124 the complete protein complement of all organisms
125 with sequenced genomes. From the results of these
126 comparisons, profiles are constructed that indicate
127 the organismal distribution of each protein's
128 homologs. Comparison of the profiles against each
129 other reveals proteins with similar phylogenetic
130 distributions, which are frequently functionally
131 linked. In one such analysis, carried out recently on
132 57 genomes, approximately 31 billion sequence
133 comparisons were made during the construction of a
134 database of phylogenetic profiles that could be
135 searched for functionally linked genes [8]. [The scale
136 of such analyses means that their results must be
137 analyzed statistically](#) to minimize the inevitable false-
138 positive linkages that arise. Nevertheless, over the
139 past year, these approaches have begun the shift in

159

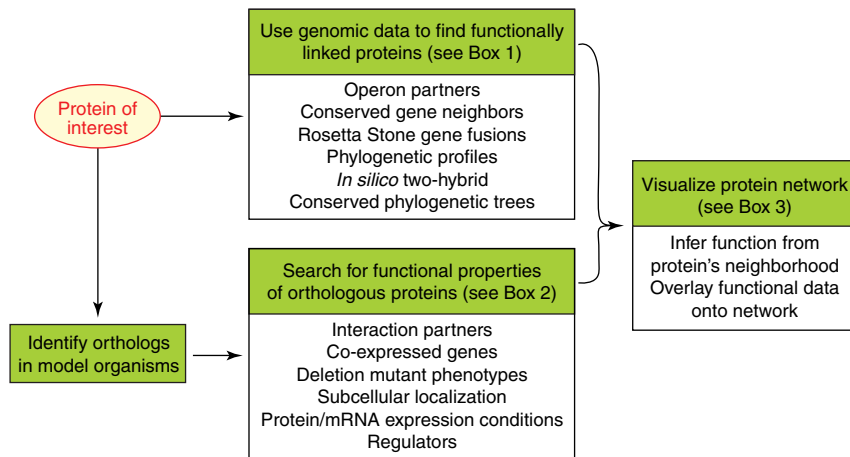
160 status from specialized research topics to publicly
161 accessible research tools. Several internet servers have
162 been created where these methods can be explored
163 and where functional linkages for a protein of
164 interest can be found, accompanied by estimates of
165 the confidence in the predictions. Several such web
166 servers are listed in Box 1.

167

168 Computational genetics approaches have proven
169 useful for several 'real-world' cases. One such case is
170 the computational identification of the archaeal
171 exosome, which was achieved using a combined
172 analysis of gene sequence homology and gene order
173 [17]. Another recent example is the discovery of
174 functional displacements of thiamin biosynthesis
175 genes [18]. In this study, candidates for gene
176 displacements in thiamin biosynthetic pathways
177 were identified using comparative genomics. Pairs of
178 genes that might substitute for each other in these
179 pathways were first identified by their anti-correlated
180 phylogenetic distributions, the involvement of these
181 genes in the biosynthetic pathway were validated
182 experimentally.

182

183 The second approach, which is poised for more
184 widespread adoption, is the 'borrowing' of function
185 from orthologs in better-characterized model
186 organisms [19], [such as the yeast *Saccharomyces
187 cerevisiae*, the nematode worm *Caenorhabditis elegans*,
188 the fly *Drosophila melanogaster*, and the bacterium
189 *Escherichia coli*](#). Given a protein of interest, it is worth
190 attempting to identify the equivalent protein in a
191 model organism and then searching for available
192 functional data among the rich functional genomics
193 and proteomics databases available for model
194 organisms. The identification of the equivalent



BioSilico

Figure 1. A flowchart describing the general 'genome-down' steps for identifying protein function computationally. Two parallel strategies exist: comparative genomics approaches for identifying linkages to other proteins, and mapping the protein of interest into an organism from which abundant functional genomics data are available and then assigning function. When using either strategy, it is often useful to examine the local network of proteins around the protein of interest, allowing the neighbors' functions to clarify the central protein's role.

protein in a model organism may not be trivial — for example, the highest-scoring BLAST match in the genome may not actually be to a protein of equivalent function — and so one typically would wish to find an orthologous, rather than simply a homologous, protein.

Orthologous genes, defined as homologous genes that are separated by speciation events [20], are typically more functionally equivalent than paralogs, defined as homologous genes separated by a duplication event. Thus, it can often be important to distinguish between these two categories of homologous genes. Even paralogs can give strong hints as to function, providing the basis for the usefulness of protein-domain databases, with the caveat that the precise functions of paralogs occasionally differ and may therefore mislead if relied upon exclusively. Identifying orthologs and paralogs requires the calculation of rooted phylogenetic trees, with outgroups, from which one can distinguish gene duplication events from speciation events. This approach is difficult to automate, and hence hard to scale to complete genomes, so several heuristic approaches have been developed that approximately identify orthologous genes.

One such heuristic approach for finding orthologs in an imperfect, yet rapid and easy, fashion has been developed by Remm and colleagues [21]. The approach is an improvement on the notion of finding 'bi-directional best hits' (BBHs). BBHs are proteins from two genomes, each of which is the top-scoring BLAST match of the other when searched in the appropriate genome [14]. Remm *et al.*'s improvement, termed InParanoid, is to recognize that many genes have been duplicated and that the

236

duplications blur the ability to identify orthologs. InParanoid therefore searches for BBHs, but then also identifies proteins from the two genomes that are as similar to the BBH proteins as the two BBH proteins are to each other. In this manner, two or more potential orthologs for a protein in the other genome may be identified within one organism

With a potential ortholog in hand, one can now search for its associated functional information. Recommended functional databases are listed in Box 2 and encompass model organism mRNA expression profiles, gene deletion phenotypes, protein subcellular localization, transcription-factor specificity, genetic interactions and protein interactions.

252 Integrating functional genomics and proteomics data

253 [Data derived from](#) DNA microarrays are one of the richest sources of information about protein function. Literally thousands of microarray datasets exist in the public domain, spawning an entire field of research in interpreting the data and distilling out functional information. Much effort in analyzing the data focuses on finding groups of genes (clusters) that have tended to co-express across a variety of experiments (reviewed by Slonim [22]).

263 This approach has been useful in suggesting functions for several uncharacterized genes (see [23,24] for examples). Without additional data, however, the results often tend to be coarse-grained and uncertain. This uncertainty is primarily caused by the inherent ambiguity in the relationships between the genes' expression patterns, which result in alternate clusterings of more-or-less equivalent quality. Furthermore, many genes may be found in a

Box 2. Functional genomics data for proteins of model organisms, which are useful for estimating the functions of orthologous proteins from other systems

Prediction of orthologs

InParanoid <http://inparanoid.cgb.ki.se>
Clusters of Orthologs (COGs) <http://www.ncbi.nlm.nih.gov/COG>

mRNA expression profiles

dbEST <http://www.ncbi.nlm.nih.gov/dbEST>
SAGEmap <http://www.ncbi.nlm.nih.gov/SAGE>
Stanford Microarray Database <http://genome-www5.stanford.edu/MicroArray/SMD>
UniGene <http://www.ncbi.nlm.nih.gov/UniGene>

Protein interaction data

Bind <http://www.bind.ca>
BRITE <http://www.genome.ad.jp/brite>
Database of Interacting Proteins <http://dip.doe-mbi.ucla.edu>
GRID <http://biodata.mshri.on.ca/grid>
MIPS <http://mips.gsf.de/proj/yeast/CYGD/db/index.html>
PIMRider (*Helicobacter pylori* protein interactions) <http://pim.hybrigenics.com>

Regulatory network data

BIOCYC/METACYC <http://biocyc.org>
GeneNet <http://www.mgs.bionet.nsc.ru/mgs/systems/genenet>
KEGG <http://www.genome.ad.jp/kegg>
Promoter Database of *S. cerevisiae* <http://cgsigma.cshl.org/jian>
RegulonDB http://www.cifn.unam.mx/Computational_Genomics/regulondb
TRANSFAC <http://transfac.gbf.de/TRANSFAC>
Yeast transcription factor targets http://web.wi.mit.edu/young/regulator_network

Model organism mutant phenotypes

Comprehensive yeast genome database <http://mips.gsf.de/genre/proj/yeast/index.jsp> **[Could you please check this link, it's not working for me.]**
FlyBase <http://flybase.bio.indiana.edu>
Saccharomyces Genome Database <http://www.yeastgenome.org>
WormBase <http://www.wormbase.org>
TRIPLES (yeast disruption phenotypes) <http://ygac.med.yale.edu/triples/triples.htm>

Protein subcellular localization and mRNA *in situ* hybridization data

TRIPLES (yeast protein localization) <http://ygac.med.yale.edu/triples/triples.htm>
Yeast green fluorescent protein (GFP) localization database <http://yeastgfp.ucsf.edu>
C. elegans mRNAs <http://nematode.laboratory.nig.ac.jp> **[Could you please check this link, it's not working for me.]**
D. melanogaster mRNAs <http://www.fruitfly.org/cgi-bin/ex/insitu.pl>
Xenopus laevis mRNAs http://www.dkfz-heidelberg.de/molecular_embryology/axelddb.htm

272
273 single cluster, complicating the precise definition of
274 their relationships. For these reasons, recent efforts to
275 interpret these data have [sought to increase](#) the
276 accuracy of the clusters. For example, Wu and
277 colleagues [25] increased accuracy by testing alternate
278 clustering methods and keeping track of those genes
279 that consistently associated together. They then
280 assigned functions to genes according to the well-
281 characterized genes that they consistently clustered
282 with. In this manner, Wu *et al.* predicted the
283 involvement of five genes in rRNA processing and
284 verified these predicted functions experimentally.
285 By integrating microarray data with other
286 functional genomics data, the quality of the
287

288
289 discovered relationships [\[between clustered](#)
290 [proteins?\]](#) has recently been improved considerably.
291
292 Statistical approaches for assigning protein
293 function from disparate sorts of data have been
294 explored in the past (e.g. [26–29]), but Troyanskaya
295 and colleagues [30] took this analysis a step further
296 by having experts in the field (mostly curators of the
297 *Saccharomyces* Genome Database) estimate the
298 accuracy of the different classes of functional
299 genomics data. Rather than directly comparing
300 different clustering results (as in [25]), Troyanskaya
301 and colleagues [30] then integrated the [functional](#)
302 [genomics data](#) according to the expert-assigned
303 accuracies using a probabilistic approach. On the

Box 3. Software for visualizing complex protein networks

General tools for network visualization

Graphlet	http://www.infosun.fmi.uni-passau.de/Graphlet/
Graphviz	http://www.research.att.com/sw/tools/graphviz/
Pajek	http://vlado.fmf.uni-lj.si/pub/networks/pajek/

Tools customized for biological networks

BioLayout	http://maine.ebi.ac.uk:8000/services/biolayout/
Cytoscape	http://www.cytoscape.org/
InterViewer3	http://wilab.inha.ac.kr/protein/
Large Graph Layout (LGL)	http://bioinformatics.icmb.utexas.edu/lgl
Osprey	http://biodata.mshri.on.ca/osprey

304
305
306
307
308
309
310
311
312
313
314
315
316

basis of these weighted data, genes were assigned to the most appropriate functional categories from the Gene Ontology project, with results available on the internet (Magic; listed in Box 1).

317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352

A second recent improvement in clustering builds on the reasonable assumption that groups of co-expressed genes are often controlled by the same set of regulatory systems. Segal and colleagues [31] exploited this [assumption](#) and took a non-obvious approach to finding gene systems: they simultaneously searched both for sets of co-expressed genes and their corresponding regulatory networks. Doubling the search problem, which at first seems to add difficulty, may actually simplify the problem by requiring mutually consistent networks and gene clusters. The algorithm works as follows: initially, the genes are simply clustered according to their expression profiles. Then, a repetitive procedure begins in which known regulatory genes are assembled into simple networks of activators and repressors that can best explain the gene expression patterns within each gene cluster. After the best set of networks has been found, the original genes are redistributed among the clusters, assigning each gene to a cluster according to how well the cluster's associated regulatory network predicts the gene's expression. Then, these two processes are alternated: first constructing the optimal regulatory network for each cluster then reassigning genes among the clusters according to the networks. The program eventually converges upon sets of co-expressed genes and their candidate regulatory networks. Unlike simply clustering the genes, this approach produces interesting, and more important, testable hypotheses that potentially explain why the genes cluster as they do.

353 Visualizing and navigating complex proteomics 354 data

355 The most intimidating aspect of working with
356 proteomics and genomics data is often the inherent

369

370 large scale and complexity [of the task](#). Nowhere is
371 this more apparent than [in attempts to unravel](#)
372 networks of proteins or genes, whose tangled sets of
373 connections are complex in the extreme. To take full
374 advantage of the data requires that a protein be
375 viewed in context, and that at least the most relevant
376 interaction partners be organized into a single
377 coherent view.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401

A useful technique for viewing networks, which was originally derived from an algorithm in computer sciences [32], has been to model proteins as an abstract objects at some location in two- or three-dimensional space. These objects are connected by springs whenever proteins are known to be linked together, and positioned so as to minimize the spring energies. It is important to realize that only the network connections represent real observations, and that the layout represents an attempt to summarize all of these perhaps conflicting associations. Therefore, the resulting distances between proteins in the network can be variable, and may even vary stochastically if the network layout is repeated. However, the tendency is for proteins that function together to be positioned close in space, whereas those that are not intimately linked tend to be further apart. In this manner, two proteins that are linked to the same set of other proteins will often be positioned adjacently, even if the two proteins are not themselves directly linked. This approach allows layouts of very complex networks, but is also famous for producing incomprehensibly complicated 'spaghetti-like' diagrams.

402
403
404
405
406
407
408
409

The difficulty of visualizing these complex datasets has led several groups to develop computer programs that allow us to visualize, and even interactively navigate, these networks more effectively. Several new network-visualization tools are listed in Box 3. Most of these new tools retain the spring-based layout approach, but use modifications to improve the visual esthetic and interpretability. Such

410 modifications include indicating different types of
411 proteins or linkages with different symbols,
412 increasing the separation between major components
413 of the network, allowing interactive navigation or
414 manipulation of the visual field of view, and even
415 collapsing proteins with more-or-less equivalent
416 interactions into single objects in the network [33] to
417 simplify the resulting network.

418 Some of the newest visualization tools, such as
419 Cytoscape and Large Graph Layout (LGL), also allow
420 the overlay of other forms of functional genomics
421 data onto the network. In this manner, protein and
422 mRNA expression levels can be simultaneously
423 viewed together with the relationships between the
424 proteins, allowing an investigator a visual summary
425 of the behavior of the system. Such a visualization
426 can allow a much more logical analysis of both the
427 expression data and the interaction data by allowing
428 the investigator to see the changes in gene expression
429 in the light of the regulatory and physical
430 interactions between the genes. [A researcher can also](#)
431 [use such tools search](#) for connected regions of the
432 interaction network that show coordinated changes
433 in gene expression patterns [34]. These tools tend to
434 be [most useful as part of a](#) genome-down approach
435 for the simple reason that a protein's final position in
436 the network map relies on maximally satisfying all of
437 the relationships in which it participates. In this
438 manner, the dominant trends in the protein's
439 relationships tend to be reinforced and suppress the
440 less-confident or less-well-observed relationships, in
441 effect providing some filtering to an otherwise very
442 complex set of relationships.

443 **Conclusions and future outlook: is the time ripe** 444 **for a central repository of protein function?**

445 The approaches discussed here provide general
446 frameworks for discovering protein function by
447 computationally integrating many distinct types of
448 data. Many types of data exist that have yet to be
449 extensively incorporated into these approaches. Such
450 datasets include protein structures and metabolite
451 and protein expression data. Protein structures
452 provide rich information about molecular aspects of
453 protein function, and it should be reasonably
454 straightforward to begin to incorporate these
455 functional inferences with those derived from
456 functional genomics and protein interaction data.
457 Little metabolite expression data exist in the public
458 domain, as useful as they would be, for example, in
459 more precisely characterizing knockout phenotypes.
460 Similarly, protein expression data are accumulating
461 rapidly in public and private laboratories, yet few of

462 this data are publicly available, curtailing the
463 development of algorithms for data analysis.

464 We expect that protein expression data will be
465 invaluable for many of the same reasons that DNA
466 microarray data are useful: they provide systematic
467 measurements of the major changes in the cell and
468 allow direct characterization of a large fraction of
469 expressed proteins. The field of functional genomics
470 has benefited tremendously from publicly accessible
471 genome sequence data and from centralized DNA
472 microarray databases (e.g. the Stanford Microarray
473 Database), and it is unfortunate that no equivalent
474 exists for proteomics. No doubt the field of
475 proteomics would profit greatly from an extensive
476 public database of protein expression data
477 contributed to by the community of proteomics
478 scientists.

479 Perhaps an equally pressing need is that for a
480 central repository of protein function data, storing
481 both experimentally determined and
482 computationally predicted functions. The biological
483 community has long had the luxury of community
484 databases that archive primary sequence, structure
485 and mRNA expression data. By contrast, the
486 distillation of functional information from these data
487 is scattered through a myriad of separate publications
488 and web servers. Several more specialized databases,
489 notably the model organism databases listed in Box 2
490 and open format sequence databases such as
491 SwissProt, have made admirable strides towards
492 cataloging this functional data, but only a small
493 fraction of computational functional analysis has
494 been included. The centralization of this
495 information, with uniformity of formats and access,
496 would open up the work of computational biologists
497 and functional genomicists to the community as a
498 whole. Most importantly, this would allow the full
499 weight of evidence for each function to be examined
500 at once. It would seem the time is ripe for
501 systematically acquired protein functions to be
502 archived systematically.

503 **Acknowledgements**

504 This work was supported by a Packard Fellowship and
505 grants from the National Science Foundation
506 (0219061,0241180), the Welch Foundation (F-1414),
507 and the Texas Advanced Research Program.

508 **References**

- 509 1 Huynen, M. *et al.* (2000) Exploitation of gene context.
510 *Curr. Opin. Struct. Biol.* 10, 366–370
- 511 2 Marcotte, E.M. (2000) Computational genetics: finding
512 protein function by nonhomology methods. *Curr. Opin.*
513 *Struct. Biol.* 10, 359–365

514 3 Marcotte, E.M. *et al.* (1999) Detecting protein function
515 and protein–protein interactions from genome sequences.
516 *Science* 285, 751–753
517 4 Enright, A.J. *et al.* (1999) Protein interaction maps for
518 complete genomes based on gene fusion events. *Nature* 402,
519 86–90
520 5 Yanai, I. *et al.* (2001) Genes linked by fusion events are
521 generally of the same functional category: a systematic
522 analysis of 30 microbial genomes. *Proc. Natl. Acad. Sci. U. S.*
523 *A.* 98, 7940–7945
524 6 Pellegrini, M. *et al.* (1999) Assigning protein functions by
525 comparative genome analysis: protein phylogenetic profiles.
526 *Proc. Natl. Acad. Sci. U. S. A.* 96, 4285–4288
527 7 Huynen, M. *et al.* (2000) Predicting protein function by
528 genomic context: quantitative evaluation and qualitative
529 inferences. *Genome Res.* 10, 1204–1210
530 8 Date, S.V. and Marcotte, E.M. (2003) Discovery of
531 uncharacterized cellular systems by genome-wide analysis of
532 functional linkages. *Nat. Biotechnol.* 21, 1055–1062
533 9 Wu, J. *et al.* (2003) Identification of functional links
534 between genes using phylogenetic profiles. *Bioinformatics* 19,
535 1524–1530
536 10 Goh, C-S. *et al.* (2000) Co-evolution of proteins with their
537 interaction partners. *J. Mol. Biol.* 299, 283–293
538 11 Pazos, F. and Valencia, A. (2001) Similarity of
539 phylogenetic trees as an indicator of protein–protein
540 interaction. *Protein Eng.* 14, 609–614
541 12 Ramani, A.K. and Marcotte, E.M. (2003) Exploiting the
542 co-evolution of interacting proteins to discover interaction
543 specificity. *J. Mol. Biol.* 327, 273–284
544 13 Dandekar, T. *et al.* (1998) Conservation of gene order: a
545 fingerprint of proteins that physically interact. *Trends*
546 *Biochem. Sci.* 23, 324–328
547 14 Overbeek, R. *et al.* (1999) The use of gene clusters to infer
548 functional coupling. *Proc. Natl. Acad. Sci. U. S. A.* 96, 2896–
549 2901
550 15 Wolf, Y.I. *et al.* (2001) Genome alignment, evolution of
551 prokaryotic genome organization, and prediction of gene
552 function using genomic context. *Genome Res.* 11, 356–372
553 16 Moreno-Hagelsieb, G. and Collado-Vides, J. (2002) A
554 powerful non-homology method for the prediction of
555 operons in prokaryotes. *Bioinformatics* 18(Suppl. 1), S329–
556 S336
557 17 Koonin, E.V. *et al.* (2001) Prediction of the archaeal
558 exosome and its connections with the proteasome and the
559 translation and transcription machineries by a comparative-
560 genomic approach. *Genome Res.* 11, 240–252
561 18 Morett, E. *et al.* (2003) Systematic discovery of analogous
562 enzymes in thiamin biosynthesis. *Nat. Biotechnol.* 21, 790–
563 795
564 19 Eisen, J.A. and Wu, M. (2002) Phylogenetic analysis and
565 gene functional predictions: phylogenomics in action. *Theor.*
566 *Popul. Biol.* 61, 481–487
567 20 Fitch, W.M. (1970) Distinguishing homologous from
568 analogous proteins. *Syst. Zool.* 19, 99–113
569 21 Remm, M. *et al.* (2001) Automatic clustering of orthologs
570 and in-paralogs from pairwise species comparisons. *J. Mol.*
571 *Biol.* 314, 1041–1052
572 22 Slonim, D.K. (2002) From patterns to pathways: gene
573 expression data analysis comes of age. *Nat. Genet.* 32(Suppl.),
574 502–508
575 23 Smith, J.J. *et al.* (2002) Transcriptome profiling to identify
576 genes involved in peroxisome assembly and function. *J. Cell*
577 *Biol.* 158, 259–271
578 24 Cheung, K.J. *et al.* (2003) A microarray-based antibiotic
579 screen identifies a regulatory role for supercoiling in the
580 osmotic stress response of *Escherichia coli*. *Genome Res.* 13,
581 206–215
582 25 Wu, L.F. *et al.* (2002) Large-scale prediction of
583 *Saccharomyces cerevisiae* gene function using overlapping
584 transcriptional clusters. *Nat. Genet.* 31, 255–265
585 26 Clare, A. and King, R.D. (2002) Machine learning of
586 functional class from phenotype data. *Bioinformatics* 18,
587 160–166
588 27 Jansen, R. *et al.* (2002) Integration of genomic datasets to
589 predict protein complexes in yeast. *J. Struct. Funct. Genomics*
590 2, 71–81
591 28 Jensen, L.J. *et al.* (2002) Prediction of human protein
592 function from post-translational modifications and
593 localization features. *J. Mol. Biol.* 319, 1257–1265
594 29 Brown, M.P.S. *et al.* (2000) Knowledge-based analysis of
595 microarray gene expression data using support vector
596 machines. *Proc. Natl. Acad. Sci. U. S. A.* 97, 262–267
597 30 Troyanskaya, O.G. *et al.* (2003) A Bayesian framework for
598 combining heterogeneous data sources for gene function
599 prediction (in *Saccharomyces cerevisiae*). *Proc. Natl. Acad. Sci.*
600 *U. S. A.* 100, 8348–8353
601 31 Segal, E. *et al.* (2003) Module networks: identifying
602 regulatory modules and their condition specific regulators
603 from gene expression data. *Nat. Genet.* 34, 166–176
604 32 Eades, P. (1984) A Heuristic for graph drawing. *Congressus*
605 *Numerantium* 42, 149–160
606 33 Ju, B.H. and Han, K. (2003) Complexity management in
607 visualizing protein interaction networks. *Bioinformatics*
608 19(Suppl. 1), I177–I179
609 34 Ideker, T. *et al.* (2002) Discovering regulatory and
610 signalling circuits in molecular interaction networks.
611 *Bioinformatics* 18(Suppl. 1), S233–S240
612
613