

# Proteomics in Genomeland

Stanley Fields

If the architect you hired to design your home brought you a blueprint that solely consisted of a long list of parts that began "windowwabeborogovestaircasedoorjub-jub....," you might start to wonder if and when you will see your new house. Some people have similar reservations about the recently "completed" human genome sequence, heralded as the "genetic blueprint" that will revolutionize biology and medicine. Deciphering how a mere  $10^7$  nucleotides results in a yeast cell—let alone how  $3 \times 10^9$  nucleotides results in Tiger Woods or Britney Spears—cannot begin until the genes have been annotated. This step includes figuring out the proteins that these genes encode and what they do for a living. But understanding how all of these proteins collaborate to carry out cellular processes is the real enterprise at hand.

In the wonderland of complete sequences, there is much that genomics cannot do, and so the future belongs to proteomics, the analysis of complete complements of proteins. Proteomics includes not only the identification and quantification of proteins, but also the determination of their localization, modifications, interactions, activities, and, ultimately, their function. Initially encompassing just two-dimensional gel electrophoresis for protein separation and identification, proteomics now refers to any procedure that characterizes large sets of proteins. The explosive growth of this field is driven by multiple forces: genomics and its revelation of more and more new proteins; powerful protein technologies, such as newly developed mass spectrometry approaches, global two-hybrid techniques, and spin-offs from DNA arrays; and innovative computational tools and methods to process, analyze, and interpret prodigious amounts of data.

The shift in thinking from genomics to proteomics comes with an appreciation of the difficulty of the task: proteins are much more complicated than nucleic acids. Unlike the decoratively challenged DNA, proteins get phosphorylated, glycosylated, acetylated, ubiquitinated, farnesylated, sulphated, linked to GPI anchors, and embellished in numerous other ways. A single

gene can encode multiple different proteins—these can be produced by alternative splicing of the mRNA transcript, by varying translation start or stop sites, or by frameshifting during which a different set of triplet codons in the mRNA is translated. All of these possibilities result in a proteome estimated to be an order of magnitude more complex than the genome. (So it may be fortunate for proteomicists that humans might have as few as six times the

number of genes that yeast have!) What is more, proteins respond to altered conditions by changing their location within the cell, getting cleaved into pieces, and adjusting their stability as well as changing what they bind to (other proteins, nucleic acids, lipids, small molecules, or other ligands). Protein levels often do not reflect mRNA levels (1), and even the presence of an open reading frame does not guarantee the existence of a protein. Lastly, a single protein may be involved in more than one process, and conversely, similar functions may be carried out by different proteins.

## Where We Are

It's worth noting that in the pre-proteomic era, thousands of proteins were exquisitely characterized—those in metabolic and signaling pathways; in the replication, transcription, and translation machinery; in secretory and cytoskeletal networks; and in a host of other cellular complexes. These functional assignments came from efforts to understand specific cellular processes, with three major factors fueling progress in the last two decades. First, a confluence of geneticists, cell biologists, biochemists, and structural biologists arrived to tackle the same problems from different directions. Second, the extraordinary conservation of fundamental cellular mechanisms

enabled insights garnered from the study of one organism to be immediately applied to all others; what's true for the walrus is true for the carpenter. Third, developments in technology—the now-standard tools of molecular biology such as DNA sequencing, recombinant DNA, and the polymerase chain reaction—sparked new experimental strategies.

With the emergence of proteomics, additional proteins are now pouring in to join those already implicated in some process or other. This extra information, however, derives not from continued small-scale analysis of a biological activity, but from much larger and more systematic studies. Proteomics, like its precursor genomics,

thus represents the emergence of a new way of doing research that is not dependent on the testing of specific models of cellular behavior. This style of science obviously does not replace, but rather will increasingly operate in tandem with, traditional biological research methods.

One general principle is that proteins prefer to hang out in the cell with others that they work with; thus, the identity of new proteins in the complexes left intact after cell lysis often provides clues to function. A big boost has come from recent advances in mass spectrometry that

allow the rapid identification of proteins separated in a two-dimensional gel or by chromatography. The mass spectrometer measures the mass of peptides (typically derived from a trypsin digestion), which is then compared to the predicted peptide masses from *in silico* digestions of sequences in genomic databases (2). Although unambiguous identification of a protein cannot always be derived from the masses of a few of its peptides, in the tandem mass spectrometer, peptide ions from the first mass spectrometer run are fragmented and identified in a second to yield the more valuable commodity of a peptide sequence. A single peptide sequence usually identifies a protein. Advances in automation, increased sensitivity, and higher throughput, combined with improved biochemical fractionations and the availability of vastly expanded databases, have extended the application of mass spectrometry to ever bigger jobs. For example, megadalton protein complexes can be purified, often with a single tagged component, and their constituents can be identified after gel electrophoresis. Such analyses have been performed on, among other complexes, the human spliceosome (3), the yeast nuclear pore complex (4), and the pea chloroplast (5). Bypassing even gel separation, the direct analysis of

The author is at the Howard Hughes Medical Institute, Departments of Genetics and Medicine, University of Washington, Seattle, WA 98195, USA. E-mail: fields@u.washington.edu.

1 protein complexes identifies components  
2 of heterogeneous protein mixtures, often  
3 using one- or two-dimensional chromatog-  
4 raphy for fractionation before analysis by  
5 mass spectrometry. Application of this  
6 procedure to a whole-cell yeast lysate  
7 identified 189 proteins (6) and more re-  
8 cently 1484 proteins, including integral  
9 membrane proteins and those of low-  
10 abundance in the cell (7).

11 The complement to mass spectrometry,  
12 the yeast two-hybrid system, has been in-  
13 creasingly “genomicized.” From its appli-  
14 cation to finding protein partners that in-  
15 teract with just one protein, the assay has  
16 been scaled up to handle, for example: 15  
17 proteins implicated in yeast mRNA splic-  
18 ing (8), 29 proteins involved in *Caeno-*  
19 *rhabditis elegans* development (9), the  
20 ~55 proteins of bacteriophage T7 (10),  
21 266 proteins of vaccinia virus (11), and  
22 even 5345 proteins of *Saccharomyces*  
23 *cerevisiae* (12). For yeast, more than 2700  
24 putative interactions involving at least  
25 2000 different proteins have been identi-  
26 fied, mostly through two-hybrid experi-  
27 ments. This set of interaction data can be  
28 visualized as protein networks, with one  
29 analysis yielding a network that encom-  
30 passes over 2300 links (13). The validity  
31 of many of the links in this network is  
32 supported by database annotations. More  
33 than 70% of characterized proteins with  
34 partners that have also been identified  
35 have been assigned a correct functional  
36 category according to the properties of  
37 these partners (compared with only 12% if  
38 the proteins in the network are kept con-  
39 stant and the links are scrambled). Thus, a  
40 protein of unknown function that binds to  
41 one of known function can be tentatively  
42 assigned to the same cellular category as  
43 its partner.

44 Protein localization within the cell can  
45 now be addressed on a genomic level. In a  
46 tour de force of transposon tagging and  
47 analysis (14), over 11,000 yeast strains  
48 were generated with more than 2,000 *S.*  
49 *cerevisiae* genes affected; indirect immu-  
50 nofluorescence was then used to deter-  
51 mine subcellular localizations for over  
52 1,300 of the tagged proteins. Biochemis-  
53 try, too, is feeling the impact of complete  
54 sequence information. The entire set of  
55 predicted yeast proteins has been fused to  
56 the “purification hook” of glutathione S-  
57 transferase (15). This set enables a bio-  
58 chemical genomics strategy in which the  
59 fusions are purified as 64 pools of 96  
60 proteins each. The pools can be assayed  
61 for any biochemical activity, and the pro-  
62 tein responsible for the activity in a pool  
63 can be quickly identified. Because the  
64 pools are derived from an array of yeast  
65 strains harboring a single gene, the gene

encoding the activity is immediately  
known.

For the expanding number of genome  
sequences available, clever algorithms  
have been developed that assign functions  
to previously unknown proteins that do not  
rely on amino acid similarity. One ap-  
proach scores the presence or absence of a  
given protein in all sequenced genomes,  
revealing sets of proteins that have co-  
evolved (that is, all members of a set are  
either present or absent in an organism),  
and are therefore likely to act in the same  
cellular process (16). A second approach is  
based on the observation that many pro-  
teins consist of two domains in one organ-  
ism, whereas the domains are two separate  
proteins in another organism (17, 18). The  
existence of the fusion, in which the two  
domains clearly interact, suggests that in  
the second organism the two separate pro-  
teins also interact. A third approach identi-  
fies cases in which multiple genomes har-  
bor the same set of neighboring genes (19,  
20), a situation implying that each set en-  
codes proteins of related function. Such  
operons in prokaryotes typically specify  
functionally linked proteins, but some ex-  
amples are also found in eukaryotes.

Although strictly speaking not a pro-  
teomics technique, DNA arrays often pro-  
vide insight into the functions of sizable  
collections of proteins. Genes that are tran-  
scriptionally co-regulated generally code  
for proteins that act in the same process, as  
demonstrated by yeast genes that operate  
in the cell division cycle (21, 22), sporu-  
lation (23), and the diauxic shift (24). Ex-  
pression profiles reveal up- or down-  
regulated mRNAs (and thus, presump-  
tively, their protein products) in disease  
processes such as cancer, and consequently  
can be used to classify tumors (25). Mi-  
croarray technology can identify classes of  
proteins—for example, membrane-bound  
and secreted proteins have been identified  
through the localization of their mRNAs  
(26), and proteins that bind to a DNA se-  
quence have been identified by their inter-  
action with a double-stranded DNA array  
(27). Microarray-based assays can also be  
used to detect polymorphisms (variations  
in the DNA), thereby associating protein  
variants with a disease state. An early ap-  
plication of this approach correctly identi-  
fied 14 of 15 patients carrying known mu-  
tations in the hereditary breast and ovarian  
cancer gene, *BRCA1* (28).

Given the current genomic and proteo-  
mic commotion, we should keep in mind  
that a protein found to be “in the  
spliceosome complex,” “interacting with  
actin,” “co-evolving with a prion protein,”  
or “upregulated in leukemia” has not been  
functionally characterized in the traditional

sense to which biologists are accustomed.  
Instead, these types of results often serve  
only to place a protein in the appropriate  
bailiwick for follow-up analysis.

### Where We're Heading

So far, most proteomic measurements have  
been performed in a cataloging mode, but  
the future will see more studies that ad-  
dress the dynamics of cellular processes.  
The protein composition of a cell is not  
static, therefore, it is crucial to obtain  
quantitative comparisons after a cell's en-  
vironment changes. Proteomic strategies  
increasingly allow such quantitative analy-  
ses to be carried out. For example, stable  
isotopes enable two protein populations to  
be labeled with either a heavy or a light af-  
finity tag, then mixed, trypsinized, and  
fractionated to enrich for subsets of pro-  
teins (29). Because the peptides in the two  
populations are identical except for the de-  
fined mass difference of the two tags,  
quantitation by mass spectrometry is pos-  
sible. These studies are in their early stages  
and their potential is tremendous. Increas-  
ingly, proteins will undergo wholesale  
analyses to probe for their various modifi-  
cations. Affinity purification approaches  
using specific antibodies, metals, lectins,  
or other reagents allow enrichment for  
modified proteins, which then can be de-  
tected by mass spectrometry (30). These  
types of strategies should make it feasible  
to follow, at the level of the proteome, a  
series of complicated cellular events such  
as those that ensue after a T cell encoun-  
ters an antigen. Advances in direct analysis  
by mass spectrometry of peptide mixtures  
generated by the digestion of complex  
protein samples will lead to an escalating  
number of protein identifications in one  
experiment. This procedure may allow  
human tissues to be used as the protein  
source and renders feasible the discovery  
of early disease markers (through the com-  
parison of the protein content of patho-  
genic cells with that of their normal coun-  
terparts).

Protein expression and purification  
technologies will continue to improve. The  
biochemical genomics strategy of purify-  
ing pools of tagged proteins will be par-  
ticularly suitable for the many bacteria that  
have had their genomes sequenced, but it  
can be applied to multicellular organisms  
as well. These and other procedures that  
make use of protein arrays will become  
commonplace. The arrays may be gener-  
ated by in vivo expression of tagged pro-  
teins, in vitro translation, peptide synthe-  
sis, or protein capture by antibodies or  
oligonucleotide aptamers. Their potential  
applications include: revealing interactions  
between proteins and between proteins and

1 small molecules (drugs) or other ligands,  
2 identifying substrates for a modifying en-  
3 zyme such as a protein kinase, and  
4 searching for enzymatic activities. A har-  
5 binger of the promise of this approach is  
6 the recent demonstration of proteins in  
7 nanoliter droplets immobilized by covalent  
8 attachment to glass slides; more than  
9 10,000 samples could be spotted onto each  
10 slide with this technique (31). The few test  
11 proteins in this array format were assayed  
12 for interactions with another protein or a  
13 small molecule, and for their phosphoryla-  
14 tion by a protein kinase. Targeted arrays  
15 will allow the identification of all of the  
16 enzymes in an organism that are able to  
17 carry out a specific modification of a sub-  
18 strate; for example, protein arrays have  
19 tested nearly the entire set of the predicted  
20 protein kinases in yeast for their activity  
21 on 17 substrates (32).

22 Protein databases will need to become  
23 much more sophisticated if they are to help  
24 scientists make sense of the staggering  
25 number of experimental measurements that  
26 will soon emerge. Demands range from  
27 tracking all of the ligands for each ana-  
28 lyzed member of a protein family (such as  
29 the SH3 domain) to cataloging all of the  
30 known substrates of each protein kinase,  
31 protein phosphatase, or other modifying  
32 activity. In addition, protein data will need  
33 to be integrated with results from expres-  
34 sion profiling, genome-wide mutation or  
35 antisense analyses, and polymorphism de-  
36 tection. As proteomic data accumulate, we  
37 will become better at triangulating from  
38 multiple disparate bits of information to  
39 gain a bearing on what a protein does in  
40 the cell. Proteomics will come of age when  
41 its revelations about formerly uncharac-  
42 terized proteins directly drive imaginative  
43 hypotheses about their functions.

#### 44 What We Need

45 For a field so laden with razzmatazz meth-  
46 ods, it is striking that the number one need  
47 in proteomics may be new technology.  
48 There are simply not enough assays that  
49 are sufficiently streamlined to allow the  
50 automation necessary to perform them on a  
51 genome's worth of proteins. Those cur-  
52 rently available barely scratch the surface  
53 of the thousands of specialized analyses  
54 biologists use every day on their favorite  
55 proteins. What we need are experimental  
56 strategies that could be termed cell bio-  
57 logical genomics, biophysical genomics,  
58 physiological genomics, and so on, to pro-  
59 vide clues to function. In addition, a pro-  
60 tein contains so many types of information  
61 that each of its properties needs to be as-  
62 sayed on a proteome-wide scale, ideally in  
63 a quantitative manner.

64 As we have argued (33), existing tech-

nology—and more importantly, the rea-  
gents (sets of genes, plasmids, strains,  
proteins, and the like) and equipment to  
handle these reagents—must rapidly  
spread from the specialized genomic and  
proteomic centers to the rest of the com-  
munity. Only when every laboratory is  
comfortable doing proteomics will its  
power be exploited fully. Moreover, the  
likelihood of new approaches increases in  
proportion to the number of investigators  
participating in the field.

An interdisciplinary spirit shall come to  
guide those excited by the global analysis  
of protein function. Geneticists need to  
talk to chemists, physiologists to physi-  
cists, cell biologists to computer scientists.  
With questions so grand, the expertise to  
answer them requires the entire spectrum  
of science. This combination of new tech-  
nology and its widespread dispersion to-  
gether with broad-ranging collaborative  
projects will culminate in the frabjous day  
when the undertaking that began with ge-  
nome sequencing reaches fruition.

#### References and Notes

1. S. P. Gygi, Y. Rochon, B. R. Franza, R. Aebersold, *Mol. Cell. Biol.* **19**, 1720 (1999).
2. W. Blackstock, in *Proteomics: A Trends Guide*, W. Blackstock, M. Mann, Eds. (Elsevier Science, London, 2000), pp. 12–17.
3. G. Neubauer *et al.*, *Nature Genet.* **20**, 46 (1998).
4. M. P. Rout *et al.*, *J. Cell Biol.* **148**, 635 (2000).
5. J. B. Peltier *et al.*, *Plant Cell* **12**, 319 (2000).
6. A. J. Link *et al.*, *Nature Biotechnol.* **17**, 676 (1999).
7. M. P. Washburn, D. Wolters, J. R. I. Yates, *Nature Biotechnol.* In press.
8. M. Fromont-Racine, J. C. Rain, P. Legrain, *Nature Genet.* **16**, 277 (1997).
9. A. J. Walhout *et al.*, *Science* **287**, 116 (2000).
10. P. L. Bartel, J. A. Roecklein, D. SenGupta, S. Fields, *Nature Genet.* **12**, 72 (1996).
11. S. McCraith, T. Holtzman, B. Moss, S. Fields, *Proc. Natl. Acad. Sci. U.S.A.* **97**, 4879 (2000).
12. P. Uetz *et al.*, *Nature* **403**, 623 (2000).
13. B. Schwikowski, P. Uetz, S. Fields, *Nature Biotechnol.*, in press.
14. P. Ross-Macdonald *et al.*, *Nature* **402**, 413 (1999).
15. M. R. Martzen *et al.*, *Science* **286**, 1153 (1999).
16. M. Pellegrini, E. M. Marcotte, M. J. Thompson, D. Eisenberg, T. O. Yeates, *Proc. Natl. Acad. Sci. U.S.A.* **96**, 4285 (1999).
17. E. M. Marcotte *et al.*, *Science* **285**, 751 (1999).
18. A. J. Enright, I. Iliopoulos, N. C. Kyripides, C. A. Ouzounis, *Nature* **402**, 86 (1999).
19. T. Dandekar, B. Snel, M. Huynen, P. Bork, *Trends Biochem. Sci.* **23**, 324 (1998).
20. R. Overbeek, M. Fonstein, M. D'Souza, G. D. Pusch, N. Maltsev, *Proc. Natl. Acad. Sci. U.S.A.* **96**, 2896 (1999).
21. R. J. Cho *et al.*, *Mol. Cell* **2**, 65 (1998).
22. P. T. Spellman *et al.*, *Mol. Biol. Cell* **9**, 3273 (1998).
23. S. Chu *et al.*, *Science* **282**, 699 (1998).
24. J. L. DeRisi, V. R. Iyer, P. O. Brown, *Science* **278**, 680 (1997).
25. C. M. Perou *et al.*, *Nature* **406**, 747 (2000).
26. M. Diehn, M. B. Eisen, D. Botstein, P. O. Brown, *Nature Genet.* **25**, 58 (2000).
27. M. L. Bulyk, E. Gentalen, D. J. Lockhart, G. M. Church, *Nature Biotechnol.* **17**, 573 (1999).
28. J. G. Hacia *et al.*, *Nature Genet.* **14**, 441 (1996).
29. S. P. Gygi *et al.*, *Nature Biotechnol.* **17**, 994 (1999).
30. O. N. Jensen, in *Proteomics: A Trends Guide*, W. Blackstock, M. Mann, Eds. (Elsevier Science, London, 2000), pp. 36–42.

31. G. MacBeath, S. L. Schreiber, *Science* **289**, 1760 (2000).
32. H. Zhu *et al.*, *Nature Genet.* **26**, 283 (2000).
33. M. Johnston, S. Fields, *Nature Genet.* **24**, 5 (2000).
34. I am especially grateful to E. Phizicky, M. Johnston, R. Aebersold, C. Manoil, and members of my laboratory for discussion and comments on the manuscript. Work from my laboratory has been supported by NIH grant P41 RR11823 from the National Center for Research Resources, and I am an investigator of the Howard Hughes Medical Institute.



1  
2 **WEB ABSTRACT**  
3 Proteomics—the large-scale analysis of  
4 a cell’s proteins—has already sup-  
5 planted genomics as the focus of bio-  
6 logical research according to Stan Fields  
7 in a lively Viewpoint. New technologies  
8 combined with traditional molecular  
9 strategies are revealing what proteins  
10 do, what they interact with and what  
11 modifications they carry. Looking to  
12 the future, even better technologies and  
13 closer collaborations between scientific  
14 disciplines will be needed to mine,  
15 analyse and compare proteomics data.  
16  
17

18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65