

**Structural Genomics Analysis:
Characteristics of atypical, typical,
and horizontally transferred folds**

Hedi Hegyi,

Jimmy Lin,

Dov Greenbaum,

&

Mark Gerstein

Department of Molecular Biophysics & Biochemistry
266 Whitney Avenue, Yale University
PO Box 208114, New Haven, CT 06520
(203) 432-6105, FAX (203) 432-5175
Mark.Gerstein@yale.edu

(Revised version, 092101)

ABSTRACT

We carried out a structural-genomics analysis of the folds and structural superfamilies in the first 20 completely sequenced genomes, focusing on the patterns of fold usage and trying to identify structural characteristics of typical and atypical folds. We assigned folds to sequences using PSI-blast, run with a systematic protocol to reduce the amount of computational overhead. On average, folds could be assigned to about a fourth of the ORFs in the genomes and about a fifth of the amino acids in the proteomes. More than 80% of all the folds in the SCOP structural classification were identified in one of the 20 organisms, with worm and *E. coli* having the largest number of distinct folds. Folds are particularly effective at comprehensively measuring levels of gene duplication, as they group together even very remote homologues. Using folds, we find the average level of duplication varies depending on the complexity of the organism, ranging from 2.4 in *M. genitalium* to 32 for the worm -- values significantly higher than those observed based purely on sequence similarity. We rank the common folds in the 20 organisms, finding that the top three are the P-loop NTP hydrolase, the ferredoxin fold, and the TIM-barrel, and discuss in detail the many factors that affect and bias these rankings. We also identify atypical folds that are "unique" to one of the organisms in our study and compare the characteristics of these folds with the most common ones. We find that common folds tend to be more multifunctional and associated with more regular, "symmetrical" structures than the unique ones. Additionally, many of the unique folds are associated with proteins involved in cell defense (e.g. toxins). We analyze specific patterns of fold occurrence in the genomes, associating some of them with instances of horizontal transfer and others with gene loss. In particular, we find three possible examples of transfer between archaea and bacteria and six between eukarya and bacteria. We make available our detailed results at <http://bioinfo.mbb.yale.edu/genome/20>.

INTRODUCTION

Structural genomics, combining the disciplines of structural biology and genomics, has emerged as a strong force in the attempt to functionally classify and annotate the genomes. It has a central concept of mapping the whole protein structure space – i.e. determining the complete protein-fold "parts list." Estimates for the total number of naturally occurring folds run somewhere between 1,000 and 10,000¹⁻³, whilst the current structural classifications divide the known structures into ~500 known folds⁴⁻⁶.

Large-scale sequence analysis of structural domains in completely sequenced microbial and eukaryotic genomes will affect both the set of proteins to be selected for experimental high-throughput structure determination, and the biological conclusions we eventually draw from the massive amount of experimental work. It is therefore timely, to perform such an analysis by comparing the sequences of the currently completed genomes to those of the already resolved and classified structural domains. Here, we survey the patterns of fold usage in the first 20 completely sequenced genomes, in the manner of a demographic census. This enables us to identify unique folds, which are potentially antibiotic targets in pathogens; shared folds, which provide information on evolutionary relatedness; common folds, which may be generic scaffolds; and overall patterns of fold usage, which may reveal aspects of protein structure and evolution beyond that found by sequence similarity. We also survey the level of gene duplication implied by the sharing of the same fold by many genes, finding, that it varies greatly between genomes.

Our work follows upon previous (mostly smaller-scale) surveys of the occurrence of folds in genomes⁷⁻¹¹, and much work on assigning folds to genomes as comprehensively as possible¹²⁻¹⁹. It also relates to a number of previous analyses in more general areas of genomics. One goal of large-scale genome analysis is to study the evolution of completely sequenced organisms by deciphering their genetic makeup through identifying orthologs and paralogs in their genomes²⁰. These studies also provide information about the conserved core of the genomes, which are necessary to the basic cellular functions of all bacteria, archaea and eukaryotes.

This survey is interesting in an evolutionary light, as it highlights those folds that are very common, as well as the unique folds within these twenty genomes. We can speculate on the evolutionary pressures, both on the structure of the folds, as well as on the functions associated with the folds to understand the make up of these populations. Moreover, we can attempt to understand what are the favorable characteristics, evolutionarily, of these folds that allow them to propagate or stay unique.

Another interesting aspect of evolution, is the relatively high frequency with which these primitive organisms incorporate foreign genes into their genomes, i.e. horizontal gene transfer²¹. These horizontally transferred genes, which provide a possible mechanism for an organism to acquire a new "part", may be represented as new folds in the organism. Or possibly, folds themselves, representing primordial self-contained proteins, may have been transferred. Analyzing a large number of closely related genomes helps to clarify this issue with greater certainty than in the past²². Large-scale genome comparison has also provided a glimpse into the evolutionary process of genome degradation in parasitic microorganisms²³.

The ultimate goal of genomics is to study biological function on a large scale. Recent success in assigning a function to a novel protein based merely on its structure, suggests that structural genomics might be useful in this endeavor. For example, Stawiski et al. identified several novel proteases based purely on their unique structural features²⁴, and Eisenstein et al. outlined a strategy to characterize 65 novel *H. influenzae* proteins through high-throughput crystallography²⁵. In terms of functional assignment, there has recently been progress based on comparing phylogenetic profiles of different gene products. These studies predict the function of an uncharacterized protein based on its consistent appearance with a protein of known function in the same genomes. Eisenberg and co-workers studied correlated evolution using phylogenetic profiles derived from 16 completely sequenced genomes, and

used these, in addition to patterns of domain fusion, to identify functionally related proteins^{26,27}. Enright *et al.* followed a similar approach and identified several unique fusion events by comparing the complete genomes of two bacteria and an archaea²⁸. Reflecting the great amount of experimental functional information available for *E. coli*, this organism's genome been studied in rather great detail in terms of functional prediction and structure-function relationships²⁹⁻³². A caveat, one has to be careful when assigning functional information through structure prediction, as promiscuous structures may often have more than one function (i.e the TIM barrel), in fact almost all superfamilies, due to local sequence variation, have multiple function.³³In addition, specific functions may be carried out by many different structures.

Finally, genomics is also driven by practical goals, such as the need to discover new antibiotics to treat emerging antibiotics-resistant bacteria. Genes that are conserved in several microbial genomes but are missing from eukaryotic genomes would be ideal targets for broad-spectrum antibiotics³⁴. Another approach is to identify species-specific genes with unique structures to reveal organism-specific biochemical pathways. Such genes are suspected to play a role in the pathogenicity of the bacteria³⁵ and could be used to develop antibiotics against specific pathogens.

Materials and Methods

Specific Databases Used in the Sequence Comparisons

Table 1A shows a list of 20 genomes we analyzed, their phylogenetic classifications, and their sizes. They represent all three domains of life (Archaea, Bacteria and Eukaryota). 19 of the 20 are single-cell organisms, and one is a eukaryote (yeast), with genome size varying from 479 (*M.genitalium*) to 6218 ORFs (yeast). The only metazoan of the twenty, *C.elegans*, has ~19000 ORFs, and the average genome size, which we denote by *G* below is 2179.

We compared the amino acid sequences of the structural domains in the SCOP classification of protein structures⁴ to the sequences of the 20 genomes. (Specifically, we used a clustered version of the SCOP database 1.39, called pdb95d, as queries. This contains 3266 distinct representative sequences, which we denote as P.) For the PSI-blast runs we also used a 90% non-redundant protein database, NRDB90³⁶, in our comparisons. This version is from December 1999 and contains 195,866 sequences (denoted as N). Both the databases (NRDB and the genome sequences) and the query sequences (SCOP domain) were masked with the SEG program using standard parameters to mask low-complexity regions^{37, 38}

Fold assignment by PSI-BLAST, Development of a Fast Hybrid Protocol

One of the goals of this work was to develop a simple, robust approach for automatically using PSI-blast³⁹ to do fold assignments for genomes in bulk.

For all our PSI-blast runs we used an inclusion threshold (*h*) of 10^{-5} , a number of iterations (*j*) of 10, and a final match threshold of 10^{-4} . These parameters are considerably more conservative than in a number of recent analyses^{12, 17, 39-41}. We were specifically concerned with guarding against false positives that would not be caught by manual checking, as we intend this to be highly automated. Furthermore, while PSI-blast, with proper masking for low-complexity regions, is known to be quite robust, the iterations occasionally run *ad absurdum* with fairly liberal parameter choices (particularly the inclusion threshold *h*) and we wished to specifically guard against this. Moreover, since we varied the size of the databases (see below) used in a variety of the runs, we wanted to try to ensure that our parameter choices resulted in significant matches in any of the databases used. We performed our PSI-blast comparisons in a number of ways:

(i) Default Protocol

We concatenated the sequences of a genome onto NRDB and used PSI-blast to run the SCOP domains as queries against them. This is the "default" way to run PSI-blast. However, it has the

drawback that every time one adds a new genome to the analysis, even a small one, one has to re-run each SCOP domain against the new genome and all of NRDB, a computationally intensive process. That is, each genome requires approximately $(N+G)PK$ pairwise comparisons, where K is the average number of iterations required by a PSI-blast comparison. (K obviously depends on many factors, including various biases both in the target database and the query, but for rough reckoning we can estimate it at $j/2 = 5$.) This is a very rough number, which we plan to use below for illustrative purposes. Using the values above it comes out to ~ 3.2 billion (3,234,074,850).

(ii) NRDB PSI-blast Profiles

We ran each SCOP query against NRDB to generate a PSI-blast profile, giving us a profile for each SCOP fold and superfamily. Then we re-ran these against the genomes without iteration, using a match threshold of 10^{-4} . (Note that because we use very conservative choices for the inclusion threshold in building up the original PSI-blast profiles, at this stage we can confidently assume that the final match threshold of 10^{-4} is selecting truly similar sequences to our original SCOP domain queries.) Note also that this is potentially a much more efficient process, since when one analyzes a new genome one only need run the profiles against each genome sequence once. That is, each new genome requires GP comparisons. (There is no K factor since there is no iteration.) Plugging in the numbers above, we get ~ 7.1 million (7,116,614).

(iii) Intra-genome Profiles

A problem with the above approach is that often the proteins that contribute most to the PSI-blast profile for a given query are in the same organism as the query. This could result, for instance, if one is searching for a protein in a family that is highly duplicated in one organism but otherwise does not have wide phylogenetic distribution. Thus, given a new genome with a highly duplicated family, one could potentially compromise sensitivity using solely NRDB generated profiles. (This would not be a problem in the default approach since one would include the genome with NRDB in the making up the of the profiles.) To get around this, while still retaining some computational efficiency for each new genome, we tried running each SCOP domain query against the genome with PSI-blast. For this protocol, for each new genome, we will require GKP comparisons, which evaluates to ~ 36 million (35,583,070) -- of course, assuming the same value for K as above, which is only approximately true.

(iv) Hybrid Protocol

For a number of select genomes, in particular *M. genitalium*, yeast and worm, we carefully compared the matches resulting from the above three protocols. We found that for the larger genomes, such as worm, use of the intra-genome profiles (protocol iii) generated quite a few additional matches beyond those found by the straight NRDB profiles (ii). In particular, using the intra-genome protocol for the worm we found 501 extra matches that were not found by the NRDB profiles (while the NRDB profiles found 576 matches that the intra-genome protocol did not find).

Combining the matches from the NRDB profiles and the intra-genome profiles (protocols ii and iii) into a new hybrid protocol resulted in essentially the same set of matches as the default PSI-blast protocol (i). For instance, for *M. genitalium*, the hybrid protocol produced at least one match for 163 different ORFs of the 483 total ORFs, whereas the default protocol produced matches for 161 different ORFs. These numbers are very similar to the values found in other PSI-blast analyses^{12, 17, 21, 40}. different ORFs. Moreover, for a new genome this was considerably more efficient than the default method, $7.1 + 3.6$ vs. 3,234 million comparisons, about 75 times more comparisons using the numbers above. To make the results of the various protocols completely clear, we make available on the web sets of matches resulting from running with the three protocols. See <http://bioinfo.mbb.yale.edu/genomes/20>. Note also that since in our hybrid protocol we are "mixing" databases for the comparisons, the precise e-values for each comparison are not exactly comparable. This is another reason for the very conservative choices we made above for our PSI-blast thresholds.

Fold assignment by FASTA, a Benchmark

As a further benchmark comparison, we ran the SCOP domains directly against the genomes using FASTA with a standard .01 e-value cutoff⁴²⁻⁴⁴. It is known that simple pairwise comparison with either FASTA or blastp is considerably less sensitive than profile-search with PSI-blast, so we did not expect this to add substantially to the number of matches that we found. However, we elected to perform the FASTA searches because for certain small compositionally biased proteins, the PSI-blast profiles may not be effective^{40,41}. Also, we felt that these would be a useful benchmark for comparison against PSI-blast. As expected, we only found a very small number of additional matches with FASTA. For instance, for the worm, the combination of the PSI-blast approaches produced at least one match for 4556 ORFs of the 19099. FASTA only added in 30 additional matches to these, considerably less than 1%, and it, of course, missed 1553 of the matches.

Tabulation in terms of SCOP Folds and Superfamilies

Using the SCOP scheme we tabulated our results in terms of distinct folds and structural superfamilies. In SCOP, for structures to have the same fold it is necessary for them to have the same overall core topology and geometric disposition of secondary structures. In contrast, a superfamily is a subset of the fold, denoting groups of proteins that have closer structural similarity and consequently probably share an evolutionary relationship⁴. We will report our specific results here separately in terms of “both SCOP folds and structural superfamilies”, henceforth known as fold.

RESULTS

Coverage of the Genome by Known Structures

Table 1A also lists the number of the ORFs in the 20 genomes that have at least one match with one of the SCOP domains, along with the ratio of these numbers and the total number of ORFs for each genome. (For a complete list of occurrences of all the folds and all the superfamilies in the 20 genomes, please see the website <http://bioinfo.mbb.yale.edu/genome/20>).

The ratio of at least partially matching ORFs varies between about 18% (for the Lyme-disease agent *B. burgdorferi*) and 34% (for *A. aeolicus* and *M. genitalium*). *M. genitalium* has often been used to benchmark the degree of fold assignment^{10, 12, 16, 40, 45}. The numbers we list for this organism are consistent with those reported in previous analyses.

Table 1A also lists the total number of amino acids in the genome "covered" by the matches and the fraction of the proteome this corresponds to (the ratio of matched and total number of amino acids). This value is surprisingly low, only about 14% for yeast and worm. Even the ‘most covered’ organisms, *A. aeolicus* and *H. influenzae*, have only slightly less than a quarter of their amino acids covered by known folds, leaving much room for either improvement in the structure prediction methods or discovery of new protein structures.

Overall Level of Duplication

The last section of Table 1A shows the level of duplication for the 20 organisms both in terms of folds (dividing the total number of domain matches by the number of different folds identified in each organism) and superfamilies (matches per superfamily). The worm has by far the highest level of fold duplication (~32), with yeast coming second with a significantly lower level, followed by *M. tuberculosis* and *E.coli*, with a fold duplication level of about 7.

Not too surprisingly, the largest number of different folds is present in the worm, followed by the most-studied microorganism, *E.coli*, while yeast is ranked only third, despite its considerably larger genome size. As for the superfamilies, *E.coli* has nearly as many as the worm (303 and 304, respectively), perhaps due to (i) a systematic bias in the structural databases, (ii) gene loss in the worm,

or (iii) folds in *E.coli* acquired by horizontal transfer from its host or other bacteria. However, the two organisms share only about two thirds (196) of their superfamilies (see the website for details).

Fold-class Specific Duplication

Table 1B also shows the total number of superfamilies and their average duplication level in the different structural classes for *A.fulgidus*, *E.coli*, yeast and worm -- representative organisms of archaea, bacteria, single-celled eukaryotes, and metazoa. One can look at this table as a subdivision of the data in Table 1 by structural class. There are clear-cut differences among the structural classes for the four organisms. In *E.coli*, the most enriched structural class is alpha/beta, while in the worm, multidomain and the small proteins are most duplicated, with a striking ~ 64X duplication level in the latter class. In yeast a similar trend can be observed, although to a lesser extent. This observation is consistent with biological observations that the majority of the small domains appear in extracellular proteins, which are required in increasing proportions to carry out the complex intercellular functions found in metazoa

There is a general depletion of the all-beta folds in the Archaea. As shown for *A.fulgidus*, only 18 superfamilies are represented, with an average duplication rate of 2.1 in this category, a relatively low value. A similar tendency can be observed in the other three archaeal genomes. Biologically, this might indicate a lesser thermostability for the all-beta structures in general, or simply reflect a lesser presence of the all-beta fold types in the last common ancestor of these organisms.

Overall Occurrence Matrix

Figure 1A shows an overview of the "occurrence matrix", the number of folds and superfamilies occurring in the six soluble fold classes for each of the 20 genomes. Each row represents a fold, each column a genome grouped by the traditional phylogenetic tree, and each cell represents the occurrence of a particular fold in a genome. The complete matrix is available in an interactive clickable form from the website. This represents the basic data from which all our fold pattern analysis is derived and provides an overall view of the structural classification used in this study. With this low-resolution diagram, although it is difficult to distinguish individual fold patterns, one can get a general sense of fold sharing among the twenty organisms.

As expected, the mixed helix and sheet classes (alpha/beta and alpha+beta) have the most universally present folds and superfamilies. The two eukaryotic genomes contain proportionately more all-alpha and all-beta folds and superfamilies than the prokaryotic ones. As previously noted, the large majority of the Small folds are present only in eukaryotes, many of them only in the metazoa worm.

Most Common Folds

Figure 1B shows a close-up of the occurrence matrix, focusing on the most frequently occurring folds and superfamilies. Two specific aspects are discussed here – the ranking biases and the top folds and superfamilies.

Factors Affecting the Ranking

In Figure 1B, to produce the ranking of the folds in terms of frequency of occurrence for the 20 genomes, we were faced with the task of arranging the folds in the occurrence matrix. There is no unique way of doing this and any method chosen introduces some form of bias. For instance, the simplest method would just order the table in terms of the raw number of matches to each fold, but these would strongly favor the large genomes, such as *C. elegans*, over the small ones, such as *M. genitalium*. Alternatively, one could rank the table purely in terms of the degree of phylogenetic conservation -- i.e. the more organisms in which a fold occurs, the higher it is in the table. However, here the ranking would be affected by the phylogenetic biases in the genomes chosen. There are many more bacterial (especially pathogen) genomes than eukaryotes. This means that folds prevalent in bacteria will tend to rank higher than those common in eukaryotes. We have developed a ranking scheme that balances a variety of factors and corrects for some obvious biases. Our scheme, described in detail in the caption to the figure,

tries to rank folds in terms of their average frequency in the main groupings of organisms (Eukaryotes, Bacteria, and Archaea), where occurrence is defined in terms of the fraction of total domains in an organism matched by a fold. (The focus on fraction of domains instead of ORFs takes into account the fact that some organisms, particularly yeast, have considerably longer ORFs than others.)

Figure 1B also shows how the highly ranked folds are connected to specific highly ranked superfamilies. When a fold is composed of many superfamilies (e.g. the TIM barrel), even if it ranks highly, the associated superfamilies may not, due to the fact that the number of folds is divided into a greater number of superfamilies. This shows how the structure of the SCOP classification itself potentially introduces a bias into the rankings. If a superfamily associated with a highly ranked fold is sufficiently different from the other members of the fold, one could potentially “split it off” and consider it as a separate fold. Doing this will decrease the ranking of the original, highly ranked fold and introduce another, lower ranking fold.

The Top-ranked Folds and Superfamilies

Based on this ranking scheme, the most abundant fold (and superfamily) in the majority of the genomes is the universally present P-loop containing NTP-hydrolase, which performs multiple biological functions. The second-ranking Ferredoxin-fold is also present in all 20 genomes; however, its most frequently occurring superfamily, 4Fe-4S Ferredoxin, is missing from several bacterial genomes. In each of the 20 genomes, at least one of the 19 superfamilies in the Ferredoxin fold is present, performing a large number of various functions, both enzymatic and non-enzymatic as explored in detail previously⁴⁶. The third-ranking fold is the TIM-barrel, also breaking down into numerous different superfamilies. This explains why even the most abundant of the TIM-barrel's superfamilies, the NAD(P)-linked oxidoreductase, ranks only 9th in the superfamily rankings. Most versatile folds defined in previous studies (TIM-barrel, Rossmann, ferredoxin, alpha-beta hydrolase, and P-loop NTP hydrolase) are all present as top folds here as well⁴⁶. Many of the most frequent folds correlated well with those identified as superfolds, i.e. folds that accommodate many distinctly different sequence families⁴⁷.

It is clear from the table that the most frequent folds and superfamilies in worm and yeast are quite different from those in the bacterial and archaeal genomes. The most abundant fold in the worm is the immunoglobulin fold, while the most abundant superfamily is the EGF/Laminin, both mostly present in extracellular, often highly repetitious proteins, providing for different functions of multicellular life.

Unique Structural Superfamilies

Table 3 shows a list of representatives for each superfamily present in only one of the 20 genomes studied here. As it appears in the table, only half of the studied organisms have unique superfamilies (we did not list the worm-specific superfamilies here; see them in a previous analysis⁴⁸). Analyzing Swissprot⁴⁹ and the non-redundant protein database (NRDB)³⁶ for the occurrence of these superfamilies helped to identify their origin: the majority of them are truly unique, occurring only in a single or a small number of organisms. The most important features of the table are summarized below:

- (a) The *B. burgdorferi*-specific outer surface protein A (ospA) was detected in seven proteins in Swissprot, all of them in this particular organism. In a sense this protein validates the idea of a unique fold in a pathogen being a drug target, as this fold is known to be the antigen for the Lyme-disease vaccine⁵⁰.
- (b) Three different domains of the enzyme copper amine oxidase are all listed as unique superfamilies occurring only in *E. coli*. (They were all detected in the same *E. coli* protein, tynA.) However, two of the three superfamilies could also be found in the human and various plant genomes (2.22.2 and 4.13.2).

- (c) The yeast metallothionein superfamily (7.38.1) was also identified in several human proteins, although not in any in the worm.
- (d) At least two superfamilies, the previously mentioned *B.burgdorferi* ospA proteins and two elastases were identified as extracellular, a rare feature for microbes.
- (e) Another unique superfamily, the flavodoxin-like cutinase, was also found only in *M. tuberculosis* in as many as seven copies. Two of them were found in neighboring ORFs, Rv3451 and Rv3452, the probable result of a recent duplication event. It is also remarkable that five out of the seven copies have an N-terminal extracellular signal. A reverse BLAST search revealed that their only homologs were found in fungi, among others in *Penicillium*, the mold that produces penicillin. One might speculate that these unique features with a unique cellular location might play an important role in the pathogenesis or the evasion of the host's immune response in this bacteria.

It appears that other pathogens might use features already 'tested' in other organisms, such as the *Pertussis* toxin, KP4 in *C. pneumoniae* or a plant pollen allergen in *M.tuberculosis*. It might be a relatively common strategy that one pathogenic microorganism could reutilize toxins that already proved to be successful in another one.

Comparing Common versus Unique Folds: typical versus atypical proteins

In our survey we found that while many of the folds and superfamilies were common (typical), some of the folds and superfamilies were unique to specific organisms. We attempted to identify possible general structural, biological, or functional explanations by comparing the common and unique folds in the survey. (We compared folds as opposed to superfamilies because of the smaller numbers.) This comparison is shown in Figure 2. We identified four main characteristics that tended to separate the common folds from the unique ones: (i) number of functions per fold (ii) nature of the function (iii) symmetry (iv) multicellularity.

(i) The common folds tend to be multifunctional. The number of functions per fold are listed in the figure⁴⁶. We can see that 19 of the 46 common folds are multifunctional. All except one of the all-beta and over half of the alpha/beta folds have more than one function, with TIM barrel having the highest number at 16 functions. However, the set of unique folds contains no multifunctional folds. Clearly, many of the common folds, may be common because they act as generic scaffolds able to carry out a variety of different functions. Of course, this always brings up an irresolvable "chicken-and-egg" issue: are the folds multifunctional because they are common or common because they are multifunctional.

(ii) Unique folds often perform specialized functions associated with cell defense. For example, we found that some of the unique folds are protease inhibitors (defensive) and toxins (offensive), as indicated in the figure. The potency and specificity of these functions logically hinges, to some degree, on the uniqueness of the fold. Toxin associated folds included the "Toxic Hairpin" (1vib, 7.002), which functions as a neurotoxin and "Yest Killer Toxin" (1kpt, chain a, 4.037). Note that two folds unique to the worm are included in this list, 1erh (7.006, topology similar to snake venom neurotoxins) and 1thw (2.019, Osmotin). While they do not clearly have functions associated with toxicity in the worm, they have toxic functions in other higher organisms not included in the figure⁵¹⁻⁵⁴. Defensive proteins include the anti-antibiotic tet-repressor (2tct, 1.094), serine protease inhibitors (1slu, 2.012 and 1pmc, 7.004) and a trypsin inhibitor (1atb, 7.022)⁵⁵⁻⁵⁸.

(iii) Many of the folds unique to only one of the twenty organisms are associated with yeast and worm, shown in bold in the table. Many of these folds are readily associated with multicellularity, as they

are only present in higher organisms. One must also take into consideration that only two eukaryotic organisms were present in this analysis.

(iv) Another characteristic among many of the common folds is that they tend to have a more symmetrical and regular structure than the unique ones. We understand that it is hard to define symmetry and "regularity" for protein folds rigorously -- this would be a subject of another project in of itself. However, simple visual examination of the unique and common folds reveals a number of obvious patterns.

Firstly, we find that the structural classes associated with common folds tend to be more regular. The alpha/beta class is the most regular, with structures required to have interleaving pattern of alpha-beta throughout. Only three of alpha/beta folds are unique whereas 19 are common. We draw all the unique and many of the 19 common alpha/beta folds in the figure. The common alpha/beta folds include well-known symmetrical structures, such as the TIM barrel and the Rossmann fold, whereas two of the three unique are clearly much more complex. The third is barnstar (1brs, 3.006), which inhibits the toxic barnase. Conversely, the structural class that is most enriched in the unique folds is that of "small" proteins. These tend to have unusual structures dominated by metal or disulfide stabilization.

Finally, we visually analyzed the all-alpha (class 1), all-beta proteins (2), and mixed non-interleaving helix-sheet proteins (4 and 5).. The common folds again are associated with a number of the well-known regular structures: the Ig fold (1ajw, 2.001), the OB fold (2prd, 2.029), the DNA-binding 3-helix bundle (1a5j, 1.004), and the 7-bladed propeller (1got, 2.051). In contrast, there are number of very complex structures associated with the unique folds. Complete images are on the associated website and we highlight a number of notable cases in the figure.

Overall Distribution of Fold Conservation

Another interesting avenue of study follows from the phylogenetic patterns of the folds, where only the presence or absence of a particular fold (or superfamily, family, etc.) in the 20 genomes is taken into consideration, and the patterns are analyzed subsequently from several viewpoints. The overall analysis of occurrence patterns is shown in Figure 3, which lists the number of superfamilies present in a given number of genomes in the six different structural classes. As expected, the alpha/beta structural class appears to be the most conserved, having 14 superfamilies common to all 20 genomes. What is more, there are only a few superfamilies in this class that appear only in one or two genomes (4 and 5, respectively). On the other hand, the all-beta, all-alpha and alpha+beta classes have many superfamilies that appear only in one or two genomes (values in these categories vary between 12 and 19, as shown in the Figure). The main reason for this, especially in the all-alpha and alpha+beta categories, is that there are many new superfamilies in these classes that appear in eukaryotes (yeast and worm here). In the Small class the large majority of the superfamilies (17) appear only in one of the 20 genomes, mostly in the worm.

A most interesting feature in this table is that the distribution in five of the six fold classes (with the exception of the Small class) does not have a "smooth tail" at the end. That is, by increasing the number of genomes, the number of conserved superfamilies does not continuously fall off; instead all have an increased value at 20 – highlighting the importance of the 38 superfamilies that are absolutely conserved throughout evolution, despite the large evolutionary diversity these 20 genomes represent. These superfamilies tend to have a disproportionately high presence in the genomes; on average about one third of all the matches in the 20 genomes belong to one of these 38 'universal' superfamilies. (However, this number varies considerably among the different genomes; in the smallest genome, *M.genitalium*, more than half the matches occurred within one of these universal superfamilies, while in *C.elegans* only about one eighth of all the matches fall into this category.) An earlier analysis we performed⁸ also indicated that many of the folds encompassing these highly conserved superfamilies tend to be superfolds⁴⁷.

Analysis of Specific Phylogenetic Patterns of Fold Occurrence

Further analysis of the overall occurrence matrix involves detailed inspection of specific patterns of fold occurrence. Some notable patterns are shown in the schematic in Figure 4. Many of these are indicative of particular evolutionary processes -- e.g. gene loss or horizontal transfer. Other patterns may indicate convergent evolution -- i.e. two folds may occur in different families of proteins that carry out the same role in different organisms but have evolved independently. Others are obvious: folds in all organisms or folds in only one. The last pattern, unique folds in certain organisms, may be useful for identifying potential drug targets. A fold present in a pathogen but not in the human genome (or in any other organism) would naturally serve as an ideal target of a highly specific drug (antibiotic or vaccine). (A detailed list of unique folds is available from the website.)

The analysis that follows shows that most of these interesting fold occurrence patterns were present in the overall occurrence matrix. The only exception is a pattern of *totally* complementary folds throughout the 20 genomes. Such a pattern is less likely to be found, as folds can be transferred between related organisms. However, we found several incomplete complementary patterns and a number of examples for horizontal fold transfer.

Gene Loss

There are a number of instances where folds (or structural superfamilies) are missing only from a single organism or clade. The most notable of these are 5 superfamilies that are missing from *Rickettsia* and present in all the other genomes.

Complementary Patterns of Fold Usage: Possible Convergent Evolution

Parts A-C of Table 2 show examples of superfamilies occurring in the different superkingdoms, performing similar or identical functions. Part A shows two superfamilies, both engaged in the control of cell division. One of them, a bacterial tubulin, is present only in archaeal and bacterial genomes (also in plants), while the other one, CKS1, a cyclin-dependent kinase, occurs only in eukaryotes.

Horizontal Transfer

It is widely recognized now that importing and reutilizing genes from foreign organisms is quite common among microbes^{59,60}. Moreover the understanding of such a process is important in attaining a clearer picture of the spread of antibiotic resistance of some bacteria. With our survey we can only suspect horizontal transfer as our results may be indicative of the biases of our survey.

Parts D and E of Table 2 list a number of possible cases of horizontal gene transfer among the three different clades. We carefully analyzed each potential candidate by collecting all proteins in Swissprot that contain domains with the same superfamily classification, and also by running reverse BLAST searches against the non-redundant (NR) protein database with the microbial ORFs as queries. Part D of the table shows 3 possible examples of such transfer from Archaea to Bacteria, while Part E lists 6 instances from Eukaryotes to Bacteria. Presently, the complexity hypothesis attempts to explain why some genes are more likely to be transferred.⁶¹ That is, there is an inverse relationship between the connectivity of the gene, i.e. the number of interactions with other genes, and the propensity to be transferred. Informational genes, those that are involved in highly organized complexes are less likely to be transferred than “operational” (i.e. housekeeping) genes. Extrapolated to folds, those folds that require other folds, either functionally, such as in the case of a fold involved in a large complex, or structurally, such as unsymmetrical folds, that need other folds to create symmetry to lower the energetic cost of the protein, are less likely to be horizontally transferred. Conversely, those folds involved in processes that do not require large complexes, such as those in our list, or folds that contain their own internal symmetry, are more likely to be transferred across organisms. With a better understanding of the evolutionary pressures upon each organisms, we may be able to deduce further meaning in the transfer of a fold from one organism to another.

Discussion

We present an analysis of 20 completely sequenced genomes in terms of their usage of protein folds. This occurrence analysis has been done very carefully, choosing the searching and iterating parameters in a way that provided a good balance between sensitivity and robustness. All our results are built upon a large table, which we call a fold occurrence matrix. Thus, we were able to rank folds in terms of their overall commonness and to broadly compare organisms in terms of sharing folds. We have also focused on specific patterns of fold usage: complementary patterns between two or more folds, unique folds in certain organisms (which are potential antibiotic targets), and horizontal transfer.

The comparison of 20 genomes in structural terms from all three kingdoms of life also provided a glimpse into the emergence and spread of new folds and superfamilies. As we noted previously⁴⁸, the worm has many specific superfamilies not present in yeast or bacteria. They are basically concerned with multicellular life, evident from the high proportion (~ 70 %) of worm-specific superfamilies that are secreted or partially extracellular. On the other hand, the eukaryote-specific superfamilies present only in the worm and yeast are typically engaged in signaling and eukaryotic-type replication, appearing mostly in multidomain proteins or protein complexes (see website for details).

The specific phylogenetic patterns reveal several interesting features of the evolution of folds and superfamilies. As it is apparent from Figure 3 and as has also been discovered by others, there is a conserved set of proteins and superfamilies that invariably are present in every genome studied so far. These completely conserved superfamilies are involved mostly in replication, and usually appear in large multidomain proteins. Furthermore, in spite of the small number of these ‘essential’ superfamilies, they amount to less than 10 percent of the total of 471 superfamilies represented in this study. However, the corresponding matches involve about one third of the total number of matching ORFs in the 20 genomes (numbers listed in Table 1). This shows that the conserved superfamilies and folds are largely over-represented in the genomes.

Another interesting point, apparent from Figure 3 is there are also many folds and superfamilies that appear in one particular or only a few organisms. We explored the 25 worm-specific superfamilies⁴⁸ and the unique superfamilies are available from the website at <http://bioinfo.mbb.yale.edu/genome/20>. Like the unique folds, many of the unique superfamilies are related to their specific life-style, e.g. the ones in *Synecocystis* are mostly related to photosynthesis, whereas pathogen bacteria often carry pathogenicity-related genes, such as the virally coded KP4 toxin in *C. pneumoniae* or the tetracycline repressor and a pollen allergen in *M.tuberculosis*. More generally, many of the unique folds and superfamilies are associated with atypical functions that in most genomes are not necessary or may be detrimental other organism. This is true almost by definition of functions associated with cell defense, since often the potency or specificity of a fold depends on its uniqueness.

Finally, we noted how the common folds tended to have a more regular and symmetrical structure than those not common. There are a number of reasons for this: (i) Symmetry is stable, economical, a low energy state and cooperative.⁶² Thus, there is higher chance that these are evolutionarily more favored. As such, it is more the rule than the exception with regard to the protein structure universe. (ii) Fold symmetry in larger proteins may be due to duplication of simpler folds through evolutionary processes, providing more complex but symmetrical folds.^{63,64} In particular, many of the current symmetric folds (e.g. the TIM barrel) could have evolved from homomultimers of simpler folds. (iii) Symmetry and regularity allow the creation of numerous but slightly different binding sites on a protein, enabling to more readily act as a generic multifunctional scaffold than one with only a single place for a site.

Future directions

Our analysis is obviously done with an incomplete list of domains, as we do not know all the protein folds. However, our analysis foreshadows the large-scale views we will have in the future after

the completion of large-scale structural genomics projects. It is worthwhile to conclude here with an enumeration of the broad types of analysis structural genomics will make possible in the future and how our work here is related to them.

(a) The complete set of protein folds will enable us to take an overall view of the occurrence of structure in nature. We will be able to see which folds occur in which organisms and which functions they are associated with. To construct the complete list of folds we will need to consider a wide variety of organisms, as it has been demonstrated that there are a number of folds specific to various phylogenetic groups.

(b) Structural genomics will much better define the actual "modules" or regions of annotation for the genome. Modules are defined by 3D structure much more precisely than by sequence patterns or motifs, and the eventual, "final" annotation of the various regions in the human genome will undoubtedly be in reference to structural modules⁸.

(c) Structural genomics will let us map the whole of protein structure space and take a global, unbiased viewpoint on the physical properties of proteins. Our view of protein structure and the conditions needed for structural stability (i.e. the size of a typical fold, the degree to which salt bridges confer thermostability, etc.) is currently strongly colored by the entries in the databanks, and this in turn is determined by the collective biases of many individual investigators following various hypothesis driven trajectories (i.e. the proteins we look at are always under the "lamppost"). It has, in fact, been shown that the proteins in the databank are NOT at all representative of those in a complete genome^{65,66}.

(d) Structural genomics will improve our understanding of distant evolution. Protein folds are among the most conserved elements in biology. In terms of folds, a great amount of redundancy and reuse occurs (as is evident in the duplication section above). Consequently, folds are ideal for probing distant evolutionary relationships, across which there is no sequence conservation. If one had a complete set of protein folds, one could see the degree to which distantly related organisms share the same underlying biochemical parts, even if the underlying genes no longer have any sequence identity.

(e) Structural genomics will enable us to see which proteins are truly generic scaffolds that occur over and over again in nature and can be used for many functions, and which are more specialised parts. In combination with gene expression and protein abundance studies^{67,68} we will also be able to see which protein folds are more highly expressed and make up the bulk of actual physical mass in a cell. Our analysis here in conjunction with other preliminary analyses suggests that the TIM barrel fold may be a most common and versatile protein part^{46,69}.

FIGURE AND TABLE CAPTIONS

Table 1A – The 20 genomes, Coverage and Duplication

Part A gives an overview of the coverage and duplication in the 20 genomes. The first column shows the 4-letter abbreviation used throughout the paper, the second column contains the full Latin names of the organisms. The literature references for the genomes are the following: Aaeo⁷⁰, Aful⁷¹, Bbur⁷², Bsub⁷³, Cpne⁷⁴, Ctra⁷⁵, Cele⁷⁶, Ecol⁷⁷, Hinf⁷⁸, Hpyl⁷⁹, Mgen⁸⁰, Mja,⁸¹, Mpne⁸², Mthe,⁸³, Mtub,⁸⁴, Phor⁸⁵, Rpro⁸⁶, Scer⁸⁷, Syne⁸⁸, Tpal⁸⁹. The third column contains the total number of ORFs in the genomes, and the fourth shows the number of ORFs that have at least one match with one of the SCOP 1.39 domains. The sixth and seventh columns show the total number of amino acids in each proteome and the number of amino acids matched by a structural domain, respectively. The fifth and eighth columns contain the percentage values of the matched ORFs and matched amino acids, respectively. (For *C. elegans*, we used the ORF file associated with it in the original publication, which contained 19099 ORFs⁷⁶ Subsequently, new versions of WormPep have come out, revising this number slightly.) The ninth and tenth columns show the number of folds and the number of superfamilies, respectively, found in the 20 genomes. The eleventh column lists the total number of matches (having eliminated the overlapping matches earlier) for each genome. The twelfth column shows the domain length for each organism. In the last two columns we calculated the fold and superfamily duplication levels, by dividing the total number of matches by the number of folds and superfamilies, respectively, present in that particular genome.

Table 1B – Represented Superfamilies and Their Average Distribution

Total number and average occurrence of the represented superfamilies in the six soluble fold classes for the genomes *A.fulgidus*, *E.coli*, yeast and worm. The last row contains the number of represented superfamilies in the 20 genomes for each class, the last column shows the total number of superfamilies in the four organisms and the total of 20 genomes.

Table 2 - Examples of interesting fold usage patterns: complementary clades and horizontal transfer.

The occurrence of dots indicates whether a particular superfamily was found in a particular genome. The table also lists the SCOP descriptions for the superfamilies, a Swissprot protein and its function containing the superfamily. **A/** Complementary clades, i.e. similar or identical functions performed by different superfamilies in the different superkingdoms between bacterial/archaeal and eukaryotic genomes. **B/** Complementary clades between bacterial and eukaryotic/archaeal genomes. **C/** Other complementary patterns, not restricted to a particular superkingdom. **D/** Examples of horizontal gene transfer between Archaea and Bacteria. **E/** Examples of horizontal gene transfer between Eukaryotes and Bacteria.

Table 3 - List of the unique superfamilies in the 20 genomes.

The unique superfamilies are shown that were manually selected after excluding FastA matches with relatively high e-values (between 0.01 and 0.001), as they appeared false positives. The first column lists the organisms where the superfamily appears. The second, third, and fourth column contains the names of the matching ORFs, the matching SCOP domains, and their classification numbers in SCOP 1.39. The next two columns contain the number of times the superfamily was found in the genome and the statistical significance of the match. The two last columns refer to the function of the superfamily.

Figure 1 – Overall fold occurrence matrix and most frequent folds and superfamilies.

The figures show two views of the "occurrence matrix" that tabulates the number of folds and superfamilies in the six soluble fold classes for each of the 20 genomes. Each row represents a fold; each column, one of the 20 genomes; and each cell represents the occurrence of a particular fold in a genome.

In both parts, the occurrence of dots indicates the presence or absence of superfamilies and folds. However, if the particular superfamily or fold is among the top ten occurrences within the genome, the cell shows a statistic relating to the matches of that fold in the genome. (Precisely, it shows $10f(i,j)$, see below.) The top occurrence in each genome is shaded in black, the next four in gray, and sixth to tenth in light gray.

The ranking scheme for folds and superfamilies is as follows: For each fold i in genome j , we first calculate the fraction of domains in the genome that have this fold: $f(i,j) = N(i,j) / D(j)$, where $N(i,j)$ is number of times fold i occurs in genome j and $D(j)$ is the estimated total number of domains in the genome. For the latter quantity we use $A(j)/170$, where $A(j)$ is the number of amino acids in the proteome of genome j (from Table 1), and 170 is an estimate of the average size of a structural domain in the PDB⁸. Notice how the calculation of $f(i,j)$ compensates for the fact that some genomes are dramatically larger than others and that the average size of a gene (in terms of amino acids and hence possible structural domains) also differs between genomes. Next, we determine an average value of $f(i)$, the fraction matched for fold i , over all genomes as follows: $f(i) = \sum_j w(j)f(i,j)$, where the weighting factor $w(j)$ is 1/6 for the two eukaryote genome, 1/12 for the four archaeal genomes, and 1/42 for the 14 bacterial genomes. The weighting factor is set so that each of the three kingdoms contributes equally to the average, and the large number of bacterial genomes does not overly skew the average. Finally, the folds or superfamilies are ranked in terms of $f(i)$.

Part A of the figure shows a schematic of the whole occurrence matrix, where the folds are first broken into major classes and then ranked in terms of $f(i)$. Part B shows a close-up of the top-ranking folds and superfamilies, including all the classes. The lines connecting the folds to the corresponding superfamilies indicate how the common folds are associated with common superfamilies. The dotted horizontal lines indicate missing lines (cuts) in the big table so that top folds in specific genomes that are not within the top total ranking can be shown. Along with each fold, the fold description and a domain identifier from SCOP 1.39⁴ are given. The entire listing is available on the website (<http://bioinfo.mbb.yale.edu/genome/20>).

Figure 2 – Comparison of Unique versus Common Folds

The two complementary tables show the members of unique and common folds in the 20 genomes. For the unique folds, the scop 1.39 id and a representative are shown in the first two columns. The third column shows the organism where the fold is unique in, as well as the number of times it occurs. The circle symbolizes toxic folds while the square symbolizes defense folds. The name is shown in the last column. For the common folds, the first two columns are the scop id and the representative. The third column shows how many functions the fold has while the last column shows the name of the fold. Percentages are shown in between the tables that indicate the composition of structural class in the set shown. On the side, different pictures of the folds are shown. Folds which are symmetric are labeled with a star; small proteins were not considered for symmetry.

Figure 3 - Distribution of the occurrence of the superfamilies among the 20 genomes.

This figure with an associated data table shows the number of SCOP superfamilies that occur in a given number of genomes. The SCOP superfamilies are divided into the usual six structural classes. For instance, the value 19 in the upper left corner of the data table denotes the 19 different all-alpha superfamilies that were found to be present in exactly one genome.

Figure 4 - Schematic.

This figure illustrates a number of interesting patterns of fold usage: (i) Present/Absent. The first pair of profiles shows two patterns in which the fold is only present in one genome, while the second pair shows patterns where the fold is absent from a single organism. The graph of the abundance of folds in each organism can be used to derive more information from the two aforementioned pairs of profiles. (ii) Complementary. The top right shows complementary patterns, in which some organisms have apparently one fold/superfamily, while other organisms have another fold/superfamily in a complementary manner. This could suggest that the two different folds/superfamilies have similar functions. However, this (complete) pattern is less likely to be found, as folds are often transferred between closely (or sometimes even remotely) related organisms. Complementary patterns in which one clade of organisms has one fold, while another one has a different fold, are more likely (middle right in the schematic). (iii) Loss/Transfer. The last two schematics show possible evidence for horizontal transfer (top of the pair) and gene loss (bottom of the pair). Horizontal transfer can be observed when one clade of organisms and just one member of the other clade have the same fold. An evolutionarily most parsimonious explanation for such a pattern is that the fold has been transferred from the dominant clade to a single organism. Gene loss can be observed when most members of the clade have the fold, whereas a few organisms do not.

1. Chothia C, Proteins — 1000 families for the molecular biologist. *Nature* 1992; 357: 543-544.
2. Govindarajan S, Recabarren R, and Goldstein RA, Estimating the total number of protein folds. *Proteins* 1999; 35: 408-14.
3. Wolf YI, Grishin NV, and Koonin EV, Estimating the number of protein folds and families from complete genome data. *J Mol Biol* 2000; 299: 897-905.
4. Murzin A, Brenner SE, Hubbard T, and Chothia C, SCOP: A Structural Classification of Proteins for the Investigation of Sequences and Structures. *J. Mol. Biol.* 1995; 247: 536-540.
5. Orengo CA, Michie AD, Jones S, Jones DT, Swindells MB, and Thornton JM, CATH--a hierarchic classification of protein domain structures. *Structure* 1997; 5: 1093-108.
6. Holm L and Sander C, Touring protein fold space with Dali/FSSP. *Nucleic Acids Res* 1998; 26: 316-9.
7. Gerstein M, A Structural Census of Genomes: Comparing Eukaryotic, Bacterial and Archaeal Genomes in terms of Protein Structure. *J. Mol. Biol.* 1997; 274: 562-576.
8. Gerstein M and Hegyi H, Comparing genomes in terms of protein structure: surveys of a finite parts list. *FEMS Microbiol Rev* 1998; 22: 277-304.
9. Frishman D and Mewes H-W, Protein structural classes in five complete genomes. *Nature Struct. Biol.* 1997; 4: 626-628.
10. Wolf YI, Brenner SE, Bash PA, and Koonin EV, Distribution of protein folds in the three superkingdoms of life. *Genome Res* 1999; 9: 17-26.
11. Lin J and Gerstein M, Whole-Genome Trees Based on the Occurrence of Folds and Orthologs: Implications for Comparing Genomes on Different Levels. *Genome Research* 2000: (in press).
12. Huynen M, Doerks T, Eisenhaber F, Orengo C, Sunyaev S, Yuan Y, and Bork P, Homology-based fold predictions for *Mycoplasma genitalium* proteins. *J Mol Biol* 1998; 280: 323-6.
13. Dubchak I, Muchnik I, and Kim SH, Assignment of folds for proteins of unknown function in three microbial genomes. *Microb Comp Genomics* 1998; 3: 171-5.
14. Fischer D and Eisenberg D, Assigning folds to the proteins encoded by the genome of *Mycoplasma genitalium*. *Proc Natl Acad Sci U S A* 1997; 94: 11929-34.
15. Fischer D and Eisenberg D, Predicting structures for genome proteins. *Curr Opin Struct Biol* 1999; 9: 208-11.
16. Muller A, MacCallum RM, and Sternberg MJ, Benchmarking PSI-BLAST in genome annotation. *J Mol Biol* 1999; 293: 1257-71.
17. Teichmann SA, Park J, and Chothia C, Structural assignments to the *Mycoplasma genitalium* proteins show extensive gene duplications and domain rearrangements. *Proc Natl Acad Sci U S A* 1998; 95: 14658-63.
18. Sanchez R and Sali A, Large-scale protein structure modeling of the *Saccharomyces cerevisiae* genome. *Proc Natl Acad Sci U S A* 1998; 95: 13597-602.
19. Salamov AA, Suwa M, Orengo CA, and Swindells MB, Genome analysis: Assigning protein coding regions to three-dimensional structures. *Protein Sci* 1999; 8: 771-7.
20. Makarova KS, Aravind L, Galperin MY, Grishin NV, Tatusov RL, Wolf YI, and Koonin EV, Comparative genomics of the Archaea (Euryarchaeota): evolution of conserved protein families, the stable core, and the variable shell. *Genome Res* 1999; 9: 608-28.
21. Wolf YI, Aravind L, and Koonin EV, Rickettsiae and Chlamydiae: evidence of horizontal gene transfer and gene exchange. *Trends Genet* 1999; 15: 173-5.
22. Doolittle RF, Feng DF, Anderson KL, and Alberro MR, A naturally occurring horizontal gene transfer from a eukaryote to a prokaryote. *J Mol Evol* 1990; 31: 383-8.
23. Andersson JO and Andersson SG, Insights into the evolutionary process of genome degradation. *Curr Opin Genet Dev* 1999; 9: 664-71.
24. Stawiski EW, Baucom AE, Lohr SC, and Gregoret LM, Predicting protein function from structure: unique structural features of proteases. *Proc Natl Acad Sci U S A* 2000; 97: 3954-8.

25. Eisenstein E, Gilliland GL, Herzberg O, Moulton J, Orban J, Poljak RJ, Banerjee L, Richardson D, and Howard AJ, Biological function made crystal clear - annotation of hypothetical proteins via structural genomics. *Curr Opin Biotechnol* 2000; 11: 25-30.
26. Pellegrini M, Marcotte EM, Thompson MJ, Eisenberg D, and Yeates TO, Assigning protein functions by comparative genome analysis: protein phylogenetic profiles. *Proc Natl Acad Sci U S A* 1999; 96: 4285-8.
27. Marcotte EM, Pellegrini M, Thompson MJ, Yeates TO, and Eisenberg D, A combined algorithm for genome-wide prediction of protein function [see comments]. *Nature* 1999; 402: 83-6.
28. Enright AJ, Iliopoulos I, Kyrpides NC, and Ouzounis CA, Protein interaction maps for complete genomes based on gene fusion events [see comments]. *Nature* 1999; 402: 86-90.
29. Zhang L, Godzik A, Skolnick J, and Fetrow JS, Functional analysis of the Escherichia coli genome for members of the alpha/beta hydrolase family. *Fold Des* 1998; 3: 535-48.
30. Tatusov RL, Mushegian AR, Bork P, Brown NP, Hayes WS, Borodovsky M, Rudd KE, and Koonin EV, Metabolism and evolution of Haemophilus influenzae deduced from a whole-genome comparison with Escherichia coli. *Curr Biol* 1996; 6: 279-91.
31. Rychlewski L, Zhang B, and Godzik A, Functional insights from structural predictions: analysis of the Escherichia coli genome. *Protein Sci* 1999; 8: 614-24.
32. Fetrow JS, Godzik A, and Skolnick J, Functional analysis of the Escherichia coli genome using the sequence- to-structure-to-function paradigm: identification of proteins exhibiting the glutaredoxin/thioredoxin disulfide oxidoreductase activity. *J Mol Biol* 1998; 282: 703-11.
33. Todd AE, Orengo CA, and Thornton JM, Evolution of function in protein superfamilies, from a structural perspective. *J Mol Biol* 2001; 307: 1113-43.
34. Galperin MY and Koonin EV, Searching for drug targets in microbial genomes. *Curr Opin Biotechnol* 1999; 10: 571-8.
35. Hacker J, Blum-Oehler G, Muhldorfer I, and Tschape H, Pathogenicity islands of virulent bacteria: structure, function and impact on microbial evolution. *Mol Microbiol* 1997; 23: 1089-97.
36. Kallberg Y and Persson B, KIND-a non-redundant protein database. *Bioinformatics* 1999; 15: 260-1.
37. Wootton JC and Federhen S, Statistics of local complexity in amino acid sequences and sequence databases. *Computers and Chemistry* 1993; 17: 149-163.
38. Wootton JC and Federhen S, Analysis of compositionally biased regions in sequence databases. *Methods Enzymol* 1996; 266: 554-71.
39. Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, and Lipman DJ, Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 1997; 25: 3389-402.
40. Teichmann SA, Chothia C, and Gerstein M, Advances in structural genomics. *Curr Opin Struct Biol* 1999; 9: 390-9.
41. Park J, Karplus K, Barrett C, Hughey R, Haussler D, Hubbard T, and Chothia C, Sequence comparisons using multiple sequences detect three times as many remote homologues as pairwise methods. *J Mol Biol* 1998; 284: 1201-10.
42. Pearson WR, Effective Protein Sequence Comparison. *Meth. Enz.* 1996; 266: 227-259.
43. Pearson WR and Lipman DJ, Improved Tools for Biological Sequence Analysis. *Proc. Natl. Acad. Sci. USA* 1988; 85: 2444-2448.
44. Pearson WR, Empirical statistical estimates for sequence similarity searches. *J Mol Biol* 1998; 276: 71-84.
45. Jones DT, GenTHREADER: an efficient and reliable protein fold recognition method for genomic sequences. *J Mol Biol* 1999; 287: 797-815.

46. Hegyi H and Gerstein M, The relationship between protein structure and function: a comprehensive survey with application to the yeast genome. *J Mol Biol* 1999; 288: 147-64.
47. Orengo CA, Jones DT, and Thornton JM, Protein superfamilies and domain superfolds. *Nature* 1994; 372: 631-4.
48. Gerstein M, Lin J, and Hegyi H, Protein Folds in the Worm Genome. *Pac. Symp. Biocomp.* 2000; 5: 30-42.
49. Bairoch A and Apweiler R, The SWISS-PROT protein sequence database and its supplement TrEMBL in 2000. *Nucleic Acids Res* 2000; 28: 45-8.
50. Li H, Dunn JJ, Luft BJ, and Lawson CL, Crystal structure of Lyme disease antigen outer surface protein A complexed with an Fab. *Proc Natl Acad Sci U S A* 1997; 94: 3584-9.
51. Barnham KJ, Dyke TR, Kem WR, and Norton RS, Structure of neurotoxin B-IV from the marine worm *Cerebratulus lacteus*: a helical hairpin cross-linked by disulphide bonding. *J Mol Biol* 1997; 268: 886-902.
52. Gu F, Khimani A, Rane SG, Flurkey WH, Bozarth RF, and Smith TJ, Structure and function of a virally encoded fungal toxin from *Ustilago maydis*: a fungal and mammalian Ca²⁺ channel inhibitor. *Structure* 1995; 3: 805-14.
53. Yun DJ, Ibeas JI, Lee H, Coca MA, Narasimhan ML, Uesono Y, Hasegawa PM, Pardo JM, and Bressan RA, Osmotin, a plant antifungal protein, subverts signal transduction to enhance fungal cell susceptibility. *Mol Cell* 1998; 1: 807-17.
54. Kieffer B, Driscoll PC, Campbell ID, Willis AC, van der Merwe PA, and Davis SJ, Three-dimensional solution structure of the extracellular region of the complement regulatory protein CD59, a new cell-surface protein domain related to snake venom neurotoxins. *Biochemistry* 1994; 33: 4471-82.
55. Grasberger BL, Clore GM, and Gronenborn AM, High-resolution structure of *Ascaris* trypsin inhibitor in solution: direct evidence for a pH-induced conformational transition in the reactive site. *Structure* 1994; 2: 669-78.
56. Brinen LS, Willett WS, Craik CS, and Fletterick RJ, X-ray structures of a designed binding site in trypsin show metal- dependent geometry. *Biochemistry* 1996; 35: 5999-6009.
57. Kisker C, Hinrichs W, Tovar K, Hillen W, and Saenger W, The complex formed between Tet repressor and tetracycline-Mg²⁺ reveals mechanism of antibiotic resistance. *J Mol Biol* 1995; 247: 260-80.
58. Mer G, Hietter H, Kellenberger C, Renatus M, Luu B, and Lefevre JF, Solution structure of PMP-C: a new fold in the group of small serine proteinase inhibitors. *J Mol Biol* 1996; 258: 158-71.
59. Pennisi E, Versatile gene uptake system found in cholera bacterium [news]. *Science* 1998; 280: 521-2.
60. Lake JA, Jain R, and Rivera MC, Mix and match in the tree of life. *Science* 1999; 283: 2027-8.
61. Jain R, Rivera MC, and Lake JA, Horizontal gene transfer among genomes: the complexity hypothesis. *Proc Natl Acad Sci U S A* 1999; 96: 3801-6.
62. Goodsell DS and Olson AJ, Structural symmetry and protein function. *Annu Rev Biophys Biomol Struct* 2000; 29: 105-53.
63. Thornton JM, Orengo CA, Todd AE, and Pearl FM, Protein folds, functions and evolution. *J Mol Biol* 1999; 293: 333-42.
64. Blundell TL and Srinivasan N, Symmetry, stability, and dynamics of multidomain and multicomponent protein systems. *Proc Natl Acad Sci U S A* 1996; 93: 14243-8.
65. Das R and Gerstein M, The Stability of Thermophilic Proteins: A Study Based on Comprehensive Genome Comparison. *Functional & Integrative Genomics* 2000; 1: 33-45.
66. Gerstein M, How representative are the known structures of the proteins in a complete genome? A comprehensive structural census. *Fold Des* 1998; 3: 497-512.

67. DeRisi JL, Iyer VR, and Brown PO, Exploring the metabolic and genetic control of gene expression on a genomic scale. *Science* 1997; 278: 680-6.
68. Gygi SP, Rochon Y, Franza BR, and Aebersold R, Correlation between protein and mRNA abundance in yeast. *Mol Cell Biol* 1999; 19: 1720-30.
69. Jansen R and Gerstein M, Analysis of the Yeast Transcriptome with Broad Structural and Functional Categories: Characterizing Highly Expressed Proteins. *Nuc. Acids Res.* 2000; 28: 1481-1488.
70. Deckert G, Warren PV, Gaasterland T, Young WG, Lenox AL, Graham DE, Overbeek R, Snead MA, Keller M, Aujay M, Huber R, Feldman RA, Short JM, Olsen GJ, and Swanson RV, The complete genome of the hyperthermophilic bacterium *Aquifex aeolicus*. *Nature* 1998; 392: 353-8.
71. Klenk HP, Clayton RA, Tomb JF, White O, Nelson KE, Ketchum KA, Dodson RJ, Gwinn M, Hickey EK, Peterson JD, Richardson DL, Kerlavage AR, Graham DE, Kyrpides NC, Fleischmann RD, Quackenbush J, Lee NH, Sutton GG, Gill S, Kirkness EF, Dougherty BA, McKenney K, Adams MD, Loftus B, Venter JC, and et al., The complete genome sequence of the hyperthermophilic, sulphate-reducing archaeon *Archaeoglobus fulgidus* [published erratum appears in *Nature* 1998 Jul 2;394(6688):101]. *Nature* 1997; 390: 364-70.
72. Fraser CM, Casjens S, Huang WM, Sutton GG, Clayton R, Lathigra R, White O, Ketchum KA, Dodson R, Hickey EK, Gwinn M, Dougherty B, Tomb JF, Fleischmann RD, Richardson D, Peterson J, Kerlavage AR, Quackenbush J, Salzberg S, Hanson M, van Vugt R, Palmer N, Adams MD, Gocayne J, Venter JC, and et al., Genomic sequence of a Lyme disease spirochaete, *Borrelia burgdorferi* [see comments]. *Nature* 1997; 390: 580-6.
73. Kunst F, Ogasawara N, Moszer I, Albertini AM, Alloni G, Azevedo V, Bertero MG, Bessieres P, Bolotin A, Borchert S, Borriss R, Boursier L, Brans A, Braun M, Brignell SC, Bron S, Brouillet S, Bruschi CV, Caldwell B, Capuano V, Carter NM, Choi SK, Codani JJ, Connerton IF, Danchin A, and et al., The complete genome sequence of the gram-positive bacterium *Bacillus subtilis* [see comments]. *Nature* 1997; 390: 249-56.
74. Kalman S, Mitchell W, Marathe R, Lammel C, Fan J, Hyman RW, Olinger L, Grimwood J, Davis RW, and Stephens RS, Comparative genomes of *Chlamydia pneumoniae* and *C. trachomatis*. *Nat Genet* 1999; 21: 385-9.
75. Read TD, Brunham RC, Shen C, Gill SR, Heidelberg JF, White O, Hickey EK, Peterson J, Utterback T, Berry K, Bass S, Linher K, Weidman J, Khouri H, Craven B, Bowman C, Dodson R, Gwinn M, Nelson W, DeBoy R, Kolonay J, McClarty G, Salzberg SL, Eisen J, and Fraser CM, Genome sequences of *Chlamydia trachomatis* MoPn and *Chlamydia pneumoniae* AR39. *Nucleic Acids Res* 2000; 28: 1397-406.
76. Consortium, Genome sequence of the nematode *C. elegans*: a platform for investigating biology. The *C. elegans* Sequencing Consortium [published errata appear in *Science* 1999 Jan 1;283(5398):35 and 1999 Mar 26;283(5410):2103 and 1999 Sep 3;285(5433):1493]. *Science* 1998; 282: 2012-8.
77. Blattner FR, Plunkett G, 3rd, Bloch CA, Perna NT, Burland V, Riley M, Collado-Vides J, Glasner JD, Rode CK, Mayhew GF, Gregor J, Davis NW, Kirkpatrick HA, Goeden MA, Rose DJ, Mau B, and Shao Y, The complete genome sequence of *Escherichia coli* K-12 [comment] [see comments]. *Science* 1997; 277: 1453-74.
78. Fleischmann RD, Adams MD, White O, Clayton RA, Kirkness EF, Kerlavage AR, Bult CJ, Tomb JF, Dougherty BA, Merrick JM, and et al., Whole-genome random sequencing and assembly of *Haemophilus influenzae* Rd [see comments]. *Science* 1995; 269: 496-512.
79. Tomb JF, White O, Kerlavage AR, Clayton RA, Sutton GG, Fleischmann RD, Ketchum KA, Klenk HP, Gill S, Dougherty BA, Nelson K, Quackenbush J, Zhou L, Kirkness EF, Peterson S, Loftus B, Richardson D, Dodson R, Khalak HG, Glodek A, McKenney K, Fitzgerald LM, Lee

- N, Adams MD, Venter JC, and et al., The complete genome sequence of the gastric pathogen *Helicobacter pylori* [see comments] [published erratum appears in Nature 1997 Sep 25;389(6649):412]. Nature 1997; 388: 539-47.
80. Fraser CM, Gocayne JD, White O, Adams MD, Clayton RA, Fleischmann RD, Bult CJ, Kerlavage AR, Sutton G, Kelley JM, and et al., The minimal gene complement of *Mycoplasma genitalium* [see comments]. Science 1995; 270: 397-403.
 81. Bult CJ, White O, Olsen GJ, Zhou L, Fleischmann RD, Sutton GG, Blake JA, FitzGerald LM, Clayton RA, Gocayne JD, Kerlavage AR, Dougherty BA, Tomb JF, Adams MD, Reich CI, Overbeek R, Kirkness EF, Weinstock KG, Merrick JM, Glodek A, Scott JL, Geoghagen NSM, and Venter JC, Complete genome sequence of the methanogenic archaeon, *Methanococcus jannaschii* [see comments]. Science 1996; 273: 1058-73.
 82. Himmelreich R, Hilbert H, Plagens H, Pirkl E, Li BC, and Herrmann R, Complete sequence analysis of the genome of the bacterium *Mycoplasma pneumoniae*. Nucleic Acids Res 1996; 24: 4420-49.
 83. Smith DR, Doucette-Stamm LA, Deloughery C, Lee H, Dubois J, Aldredge T, Bashirzadeh R, Blakely D, Cook R, Gilbert K, Harrison D, Hoang L, Keagle P, Lumm W, Pothier B, Qiu D, Spadafora R, Vicaire R, Wang Y, Wierzbowski J, Gibson R, Jiwani N, Caruso A, Bush D, Reeve JN, and et al., Complete genome sequence of *Methanobacterium thermoautotrophicum* deltaH: functional analysis and comparative genomics. J Bacteriol 1997; 179: 7135-55.
 84. Cole ST, Brosch R, Parkhill J, Garnier T, Churcher C, Harris D, Gordon SV, Eiglmeier K, Gas S, Barry CE, 3rd, Tekaiia F, Badcock K, Basham D, Brown D, Chillingworth T, Connor R, Davies R, Devlin K, Feltwell T, Gentles S, Hamlin N, Holroyd S, Hornsby T, Jagels K, Barrell BG, and et al., Deciphering the biology of *Mycobacterium tuberculosis* from the complete genome sequence [see comments] [published erratum appears in Nature 1998 Nov 12;396(6707):190]. Nature 1998; 393: 537-44.
 85. Kawarabayasi Y, Sawada M, Horikawa H, Haikawa Y, Hino Y, Yamamoto S, Sekine M, Baba S, Kosugi H, Hosoyama A, Nagai Y, Sakai M, Ogura K, Otsuka R, Nakazawa H, Takamiya M, Ohfuku Y, Funahashi T, Tanaka T, Kudoh Y, Yamazaki J, Kushida N, Oguchi A, Aoki K, and Kikuchi H, Complete sequence and gene organization of the genome of a hyper-thermophilic archaeobacterium, *Pyrococcus horikoshii* OT3 (supplement). DNA Res 1998; 5: 147-55.
 86. Andersson SG, Zomorodipour A, Andersson JO, Sicheritz-Ponten T, Alsmark UC, Podowski RM, Naslund AK, Eriksson AS, Winkler HH, and Kurland CG, The genome sequence of *Rickettsia prowazekii* and the origin of mitochondria [see comments]. Nature 1998; 396: 133-40.
 87. Goffeau A, Barrell BG, Bussey H, Davis RW, Dujon B, Feldmann H, Galibert F, Hoheisel JD, Jacq C, Johnston M, Louis EJ, Mewes HW, Murakami Y, Philippsen P, Tettelin H, and Oliver SG, Life with 6000 genes [see comments]. Science 1996; 274: 546, 563-7.
 88. Kaneko T, Sato S, Kotani H, Tanaka A, Asamizu E, Nakamura Y, Miyajima N, Hirosawa M, Sugiura M, Sasamoto S, Kimura T, Hosouchi T, Matsuno A, Muraki A, Nakazaki N, Naruo K, Okumura S, Shimpo S, Takeuchi C, Wada T, Watanabe A, Yamada M, Yasuda M, and Tabata S, Sequence analysis of the genome of the unicellular cyanobacterium *Synechocystis* sp. strain PCC6803. II. Sequence determination of the entire genome and assignment of potential protein-coding regions. DNA Res 1996; 3: 109-36.
 89. Fraser CM, Norris SJ, Weinstock GM, White O, Sutton GG, Dodson R, Gwinn M, Hickey EK, Clayton R, Ketchum KA, Sodergren E, Hardham JM, McLeod MP, Salzberg S, Peterson J, Khalak H, Richardson D, Howell JK, Chidambaram M, Utterback T, McDonald L, Artiach P, Bowman C, Cotton MD, Venter JC, and et al., Complete genome sequence of *Treponema pallidum*, the syphilis spirochete [see comments]. Science 1998; 281: 375-88.