# NEWS & VIEWS

# Gastrogenomic delights: A movable feast

**The complete genome sequences of *Escherichia coli* and *Helicobacter pylori* provide insights into the biology of these species.**

JONATHAN A. EISEN[1], DALE KAISER[2] & RICHARD M. MYERS[3]

RECENTLY, WE BIOLOGISTS have been treated to a feast of the complete genome sequences of two gut bacteria: *Helicobacter pylori* reported by Tomb *et al.*[1] in *Nature* and *Escherichia coli* reported by Blattner *et al.*[2] in *Science*. Complete sequences of eight microbes have now been published (Table 1), and there are over 30 additional projects underway and slated for completion in the next 12–18 months. The finished genome sequence of *E. coli* — metabolic generalist, workhorse of biochemical genetics, molecular biology and biotechnology, and occasional pathogen — has special, almost emotional, significance to today's biologists, many of whom have grown up with its cultures in one form or another. By contrast, *H. pylori* — metabolic specialist, gastric pathogen and causative agent of peptic ulcers — is a relative newcomer to the scientific scene (Fig. 1).

There are numerous reasons for going to the trouble of determining complete and accurate genome sequences of micro-organisms. In those microbes with pathogenic properties, the total set of instructions provides a potentially powerful basis for developing vaccines and other therapeutic agents. Genomic sequences offer insights into the range of functions an organism possesses, the relative importance natural selection attaches to each function, and the organism's evolutionary history. In addition, the availability of complete genome sequences has spawned an enormous array of creative approaches for global functional analysis of genes and gene networks. There is particular virtue in having contiguous sequences of an entire genome; not only is it possible to predict all or almost all of the proteins that are present in the organism, but what is absent also becomes meaningful.

The genome sequence of an organism is like the Rosetta stone: it is impressive to see, but it must be translated to have value. The most important initial steps in translating a genome are identifying all of the genes and assigning functions to them. Genes can be identified by genetic and biochemical experiments or predicted by computational analysis of the genome sequence. Functions of genes can be assigned also by experi-

mental and computational methods, but accurate prediction of function based solely on sequence information is not so straightforward. In the case of *E. coli*, computational prediction of gene function is less important because of the vast wealth of genetic and biochemical data collected from this organism over the last fifty years[3]. However, for *H. pylori* and for most of the species for which complete genome sequences are published, far less experimentally derived functional information is available. Thus, analysis of these genomes, and most of those that will be sequenced in the future, depends heavily on computational methods.

Tomb *et al.* use the BLAZE program[4] to assign a function to each predicted *H. pylori* gene based on the function of a previously characterized gene in the sequence database that is most similar in sequence to the predicted gene, but only if the likelihood of the match is much higher than that expected by chance. Blattner's group go one step further. They identify multiple similar sequences in existing databases, and if most of these genes appear to have the same physiological role, this function is assigned to the new gene. If the top scoring sequences have differ-

ent physiological roles, attempts are made to identify a common denominator, such as transport activity, and this general activity, with unknown specificity, is then assigned to the new gene. Although both approaches are likely to result in correct functional assignments for most genes, there are many cases in which either approach will lead to incorrect predictions.

One example where caution seems warranted is in the prediction that *H. pylori* is capable of mismatch repair, based on the assignment of methyl transferase, MutS and UvrD functions to several of its genes[1]. However, it is unlikely that this DNA repair process is present in *H. pylori* because its genome sequence does not contain a homologue of MutL, a protein required for mismatch repair in all organisms studied from bacteria to humans[5]. Furthermore, phylogenetic analysis suggests that there has been an ancient duplication in the MutS gene family, and that the "MutS" gene (HP0621) in *H. pylori* is not an orthologue (a gene originating from a speciation event), but is rather a paralogue (a gene originating from a gene duplication event) of the *E. coli* MutS gene (Fig. 2). Genes that are orthologues of the *E. coli* MutS gene (Fig. 2, blue) are absolutely required for mismatch repair in many bacterial species. By contrast, the MutS paralogues in bacteria (Fig. 2, red) have
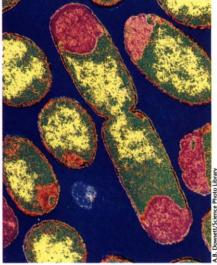


**Fig. 1**  Transmission electron micrographs of *Helicobacter pylori* (left) and *Escherichia coli* (right).

| Table 1 | Complete Genomes | | | | |
|---|---|---|---|---|---|
| Species | Classification | Size (Mb) | ORFs | Ref. | Web Site |
| **Bacteria** | | | | | |
| Mycoplasma genitalium | LowGC gram positive | 0.58 | 470 | Fraser et al.[9] 1995 | http://www.tigr.org/tdb/mdb/mgdb/mgdb.html |
| Mycoplasma pneumoniae | LowGC gram positive | 0.82 | 679 | Himmelreich et al.[10] 1996 | http://www.zmbh.uni-heidelberg.de/M_pneumoniae/MP_Home.html |
| Escherichia coli | Proteobacteria (γ) | 4.60 | 4288 | Blattner et al.[2] 1997 | http://www.genetics.wisc.edu:80/index.html |
| Haemophilus influenzae | Proteobacteria (γ) | 1.83 | 1743 | Fleischman et al.[11] 1995 | http://www.tigr.org/tdb/mdb/hidb/hidb.html |
| Helicobacter pylori | Proteobacteria (ε) | 1.67 | 1590 | Tomb et al.[1] 1997 | http://www.tigr.org/tdb/mdb/hpdb/hpdb.html |
| Synechocystis sp. | Cyanobacteria | 3.57 | 3168 | Kaneko et al.[12] 1996 | http://www.kazusa.or.jp/cyano/cyano.html |
| **Archaea** | | | | | |
| Methanococcus jannascii | Euryarchaeota | 1.66 | 1738 | Bult et al.[13] 1996 | http://www.tigr.org/tdb/mdb/mjdb/mjdb.html |
| **Eukaryote** | | | | | |
| Saccharomyces cerevisiae | Fungi | 12.07 | 5885 | Goffeau et al.[14] 1997 | http://genome-www.stanford.edu/Saccharomyces/ |

no known function. Why was the HP0621 gene of *H. pylori* called MutS, and not identified as a MutS paralogue? Analysis of the database search used by Tomb's group (see their web site, Table 1) indicates that the gene was given this designation because its highest sequence similarity was with the gene sll1772 from *Synechocystis* sp. (strain PCC6803), a cyanobacterium. The researchers annotating the *Synechocystis* sp. genome sequence earlier gave the name MutS to gene sll1772, because it scored as highly similar to MutS in a similar type of analysis. However, gene sll1772 is only one of two MutS-like genes in *Synechocystis* sp. A second gene (gene sll1165) predicted from its genome sequence is much more similar to MutS of *E. coli*, and is the likely MutS orthologue in *Synechocystis* sp. *H. pylori*, for unknown reasons, does not encode an orthologue of the *MutS* genes known to be involved in mismatch re-

pair. As this example shows, database errors are often self-propagating.

The difficulty in assigning function on the basis of sequence data is likely to be widespread, particularly because so many microbial genome sequences are forthcoming. Some simple precautions may help to alleviate the problem. Perhaps the most obvious rule is to avoid assuming that a function assigned to a sequence is correct just because it already appears in a database. A second simple precaution is to recognize that sequence similarity indicates only the *potential* for a biochemical activity. Close similarity does not readily identify the physiological role for a protein and is not definitive evidence that two proteins have the same biochemical activity. Likewise, the absence of a homologous gene in a whole genome sequence does not necessarily mean that the activity is absent in the organism.

The MutS story and other examples provide evidence that classifying members of multigene families is one of the most difficult parts of assigning function. Molecular phylogenetics is probably a better method for dividing multigene families into groups of orthologous genes rather than simply relying on database searches. As orthologues frequently have functions distinct from paralogues, a "phylogenomic" method is likely to improve the accuracy of function assignment to members of multigene families identified in complete genome sequences. In addition, assignment of function on the basis of DNA sequence data should become more accurate as we learn how to integrate knowledge about biochemical pathways and regulatory networks into the computational methods.

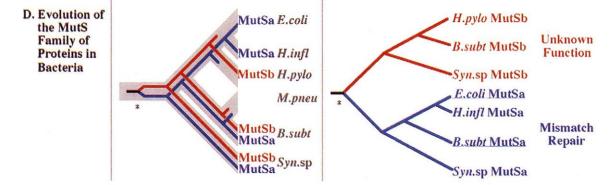In addition to stimulating predictions of the functions of individual genes, the



**Fig. 2** Reconstruction of the evolution of MutS-like proteins in bacteria using molecular phylogenetics. MutS-like protein sequences were aligned and a tree of these sequences was generated using molecular phylogenetic methods *(details are available from the authors on request)*. *left*, The tree of the MutS-like proteins (thin lines) is shown embedded within the species tree (thick grey lines). The gene duplication event (marked by an asterisk) occurred prior to the divergence of these bacterial species and led to the presence of two paralogous MutS-like subgroups (distinguished by different colors and gene subscripts a or b). Gene loss in some lineages is indicated when the MutS tree stops within the species tree. *right*, The MutS tree is extracted from the species tree and untwisted to show better the relationships among the different MutS forms. Only one lineage (labeled in blue) includes genes with established roles in mismatch repair. The genes in the second lineage (in red) have no known function. Because the *H. pylori* gene is a member of this second lineage, and because paralogous genes frequently have different functions, it should not be assigned the MutS function.

complete genome sequences of *H. pylori* and *E. coli* provide clues about their global metabolic capabilities. One striking difference between these two organisms is that *H. pylori* has far fewer genes than *E. coli*. Three other bacteria with complete sequences published also have small genomes. How can this phenomenon be explained? One argument is that organisms with broad ecological niches need more genes[6]. For example, *E. coli*, with a genome of 4.6 million base pairs, can be thought of as a metabolic generalist because it is capable of growing under a variety of conditions. It is equipped to grow in the lower gut of animals, where it meets a variety of sugars that have not been absorbed by its host's digestive tract. The lower gut is also anaerobic; *E. coli* is a facultative anaerobe, capable of fermentative metabolism. *E. coli* survives when it is released to the environment, where it can be disseminated to new hosts. It grows faster in air than in the gut, metabolizing carbon completely to carbon dioxide. Its metabolic generalism is reflected in its genome; there are many different transport proteins to accumulate dilute substrates from the gut contents. There are 700 known gene products for central intermediary metabolism, degradation of small molecules and energy metabolism.

Helping *E. coli* adjust to a variety of growth conditions are the 400 regulatory genes (some known on the basis of experiments and some attributed for reasons of sequence similarity) that constitute 4.5% of its total genome. By contrast, *H. pylori*, with a genome of only 1.66 million base pairs, is an ecological specialist, apparently living nowhere but in the mucosa of the stomach. Consistent with this restricted ecological niche, the genome sequence of *H. pylori* indicates that it is much more limited than *E. coli* in its metabolic capabilities and its regulatory networks. The genome sequence also provides clues as to how *H. pylori* survives in the highly acidic environment of the stomach. The proteins encoded by the *H. pylori* genome have twice the number of basic amino acids compared to proteins of other microbes; this may help in establishing a positive inside membrane potential. These comparisons provide but one of many valuable insights that can be gleaned from sequences of complete bacterial genomes.

We have every reason to be delighted by the feast that has just been served to us. These complete genomic sequences have a major impact on the study of these two gut bacteria, and are likely to speed up our understanding of the mechanisms by which they cause disease. Because these and the other available bacterial sequences are from widely divergent microbes, we are already getting an idea of which genes are universal and perhaps form the core of a micro-organism[7]. By contrast, as complete genome sequences from closely related pairs of microbes become available, we will learn more about mutation and recombination processes, as well as features such as codon usage, genome structure, and horizontal gene transfer, that change on a shorter evolutionary time scale[8]. Today's feast may well seem meager in comparison to the lavish smorgasbord expected in the future.

1. Tomb, J-F. *et al.* The complete genome sequence of the gastric pathogen *Helicobacter pylori*. *Nature* **388**, 539–547 (1997).
2. Blattner, F. R. *et al.* The complete genome sequence of *Escherichia coli* K-12. *Science* **277**, 1453–1462 (1997).
3. Riley, M. Functions of the gene products of *Escherichia coli*. *Micro. Rev.* **57**, 862–952 (1993).
4. Brutlag, D. L. *et al.* BLAZE: An implementation of the Smith-Waterman comparison algorithm on a massively parallel computer. *Compu. and Chem.* **17**, 203–207 (1993).
5. Modrich, P. & Lahue, R. Mismatch repair in replication fidelity, genetic recombination, and cancer biology. *Ann. Rev. Biochem.* **65**, 101–133 (1996).
6. Hinegardner, R. Evolution of genome size. In *Molecular Evolution* (ed. Ayala, F.J.) 179–199 (Sinauer Sunderland, MA, 1976).
7. Mushegian, A. R. & Koonin, E. V. A minimal gene set for cellular life derived by comparison of complete bacterial genomes. *Proc. Natl. Acad. Sci. U. S. A.* **93**, 10268–10273 (1996).
8. Lawrence, J. G. & Ochman, H. Amelioration of bacterial genomes: rates of change and exchange. *J. Mol. Evol.* **44**, 383–397 (1997).
9. Fraser, C. M. *et al.* The minimal gene complement of *Mycoplasma genitalium*. *Science* **270**, 397–403 (1995).
10. Himmelreich, R. *et al.* Complete sequence analysis of the genome of the bacterium *Mycoplasma pneumoniae*. *Nucleic Acids Res.* **24**, 4420–4449 (1996).
11. Fleischmann, R. D. *et al.* Whole-genome random sequencing and assembly of *Haemophilus influenzae* Rd. *Science* **269**, 496–512 (1995).
12. Kaneko, T. *et al.* 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.* **3**, 109–136 (1996).
13. Bult, C. J. *et al.* Complete genome sequence of the methanogenic archaeon, *Methanococcus jannaschii*. *Science* **273**, 1058–1073 (1996).
14. Goffeau, A. *et al.* The yeast genome directory. *Nature* **387** (Suppl.) 5–105 (1997).

*Departments of* [1]*Biological Sciences*
[2]*Biochemistry and* [3]*Genetics*
[1]*Stanford University and*
[2,3]*Stanford University School of Medicine*
*Stanford, California 94305, USA*

# A coat of many complements

**A novel mechanism used by pathogenic mycobacteria to enter macrophages (in low complement environments) may provide a new therapeutic target for the treatment of tuberculosis.**

YONA ZAFFRAN & JERROLD J. ELLNER

THE CAUSATIVE ORGANISM of tuberculosis, *Mycobacterium tuberculosis* (MTB), is an intracellular pathogen that resides in the alveolar macrophages of infected individuals. To initiate infection and for successful replication mycobacteria must be recognized and engulfed by macrophages. Recognition may be through several different opsonic[1] or nonopsonic pathways (see figure). Opsonins include complement protein C3b and its degradation product C3bi, which bind macrophage receptors CR1 and CR3/CR4, respectively. A recent *Science* paper by Schorey *et al.*[2] reports that the complement cleavage product C2a is used by three pathogenic strains of mycobacteria — MTB, *M. avium* and *M. leprae* — to produce opsonins that enhance invasion of macrophages. However, this mechanism is not used by rapidly growing nonpathogenic mycobacteria or nonmycobacterial intracellular pathogens such as *Leishmania mexicana* or *Listeria monocytogenes*. These findings suggest that C2a may constitute a virulence mechanism specific for pathogenic strains of mycobacteria.

C2a is generated by the classical pathway of complement activation and normally interacts with C4b to form an enzyme (convertase) that cleaves C3. Schorey and colleagues show that, in the absence of C4b, C2a can form a C3 convertase that cleaves C3 to the opsonin C3b, which is then deposited on the surface of pathogenic mycobacteria. The nature of the mycobacterial surface receptor that binds C2a is of considerable interest and may be a future therapeutic target for tuberculosis drug development. The tiny concentrations of C2a required (1–10 nM) and the potential for contribution of complement components produced by macrophages[3] may favor mycobacterial uptake in low opsonin environments such as in the lung. In addition, tuberculosis is characterized by local inflammation that is associated with systemic complement activation[4].

One possible advantage of this novel