

LETTERS

A phylogeny-driven genomic encyclopaedia of Bacteria and Archaea

Dongying Wu^{1,2}, Philip Hugenholtz¹, Konstantinos Mavromatis¹, Rüdiger Pukall³, Eileen Dalin¹, Natalia N. Ivanova¹, Victor Kunin¹, Lynne Goodwin⁴, Martin Wu⁵, Brian J. Tindall³, Sean D. Hooper¹, Amrita Pati¹, Athanasios Lykidis¹, Stefan Spring³, Iain J. Anderson¹, Patrik D'haeseleer^{1,6}, Adam Zemla⁶, Mitchell Singer², Alla Lapidus¹, Matt Nolan¹, Alex Copeland¹, Cliff Han⁴, Feng Chen¹, Jan-Fang Cheng¹, Susan Lucas¹, Cheryl Kerfeld¹, Elke Lang³, Sabine Gronow³, Patrick Chain^{1,4}, David Bruce⁴, Edward M. Rubin¹, Nikos C. Kyrpides¹, Hans-Peter Klenk³ & Jonathan A. Eisen^{1,2}

Sequencing of bacterial and archaeal genomes has revolutionized our understanding of the many roles played by microorganisms¹. There are now nearly 1,000 completed bacterial and archaeal genomes available², most of which were chosen for sequencing on the basis of their physiology. As a result, the perspective provided by the currently available genomes is limited by a highly biased phylogenetic distribution^{3–5}. To explore the value added by choosing microbial genomes for sequencing on the basis of their evolutionary relationships, we have sequenced and analysed the genomes of 56 culturable species of Bacteria and Archaea selected to maximize phylogenetic coverage. Analysis of these genomes demonstrated pronounced benefits (compared to an equivalent set of genomes randomly selected from the existing database) in diverse areas including the reconstruction of phylogenetic history, the discovery of new protein families and biological properties, and the prediction of functions for known genes from other organisms. Our results strongly support the need for systematic 'phylogenomic' efforts to compile a phylogeny-driven 'Genomic Encyclopedia of Bacteria and Archaea' in order to derive maximum knowledge from existing microbial genome data as well as from genome sequences to come.

Since the publication of the first complete bacterial genome, sequencing of the microbial world has accelerated beyond expectations. The inventory of bacterial and archaeal isolates with complete or draft sequences is approaching the two thousand mark². Most of these genome sequences are the product of studies in which one or a few isolates were targeted because of an interest in a specific characteristic of the organism. Although large-scale multi-isolate genome sequencing studies have been performed, they have tended to be focused on particular habitats or on the relatives of specific organisms. This overall lack of broad phylogenetic considerations in the selection of microbial genomes for sequencing, combined with a cultivation bottleneck⁶, has led to a strongly biased representation of recognized microbial phylogenetic diversity^{3–5}. Although some projects have attempted to correct this (for example, see ref. 5), they have all been small in scope. To evaluate the potential benefits of a more systematic effort, we embarked on a pilot project to sequence approximately 100 genomes selected solely for their phylogenetic novelty: the 'Genomic Encyclopedia of Bacteria and Archaea' (GEBA).

Organisms were selected on the basis of their position in a phylogenetic tree of small subunit (SSU) ribosomal RNA, the best sampled

gene from across the tree of life⁷. Working from the root to the tips of the tree, we identified the most divergent lineages that lacked representatives with sequenced genomes (completed or in progress)⁸ and for which a species has been formally described⁹ and a type strain designated and deposited in a publicly accessible culture collection¹⁰. From hundreds of candidates, 200 type strains were selected both to obtain broad coverage across Bacteria and Archaea and to perform in-depth sampling of a single phylum. The Gram-positive bacterial phylum *Actinobacteria* was chosen for the latter purpose because of the availability of many phylogenetically and phenotypically diverse cultured strains, and because it had the lowest percentage of sequenced isolates of any phylum (1% versus an average of 2.3%)¹¹. Of the 200 targeted isolates, 159 were designated as 'high' priority primarily on the basis of phylum-level novelty and the ability to obtain microgram quantities of high quality DNA. The genomes of these 159 are being sequenced, assembled, annotated (including recommended metadata¹²) and finished, and relevant data are being released through a dedicated Integrated Microbial Genomes database portal¹³ and deposited into GenBank. Currently, data from 106 genomes (62 of which are finished) are available.

To assess the ramifications of this tree-based selection of organisms, we focused our analyses on the first 56 genomes for which the shotgun phase of sequencing was completed. The 53 bacteria and 3 archaea (Supplementary Table 1) represent both a broad sampling of bacterial diversity and a deeper sampling of the phylum *Actinobacteria* (26 GEBA genomes). An initial question we addressed was whether selection on the basis of phylogenetic novelty of SSU rRNA genes reliably identifies genomes that are phylogenetically novel on the basis of other criteria. This question arises because it is known that single genes, even SSU rRNA genes, do not perfectly predict genome-wide phylogenetic patterns^{14,15}. To investigate this, we created a 'genome tree' (ref. 16) of completed bacterial genomes (Fig. 1) and then measured the relative contribution of the GEBA project using the phylogenetic diversity metric¹⁷. We found that the 53 GEBA bacteria accounted for 2.8–4.4 times more phylogenetic diversity than randomly sampled subsets of 53 non-GEBA bacterial genomes. A similar degree of improvement in phylogenetic diversity was seen for the more intensively sampled actinobacteria (Table 1). These analyses indicate that although SSU rRNA genes are not a perfect indicator of organismal evolution, their phylogenetic relationships are a sound predictor of phylogenetic novelty within the universal gene core present in bacterial genomes.

¹DOE Joint Genome Institute, Walnut Creek, California 94598, USA. ²University of California, Davis, California 95616, USA. ³DMSZ, German Collection of Microorganisms and Cell Cultures, 38124 Braunschweig, Germany. ⁴DOE Joint Genome Institute-Los Alamos National Laboratory, Los Alamos, California 87545, USA. ⁵University of Virginia, Charlottesville, Virginia 22904, USA. ⁶Lawrence Livermore National Laboratory, Livermore, California 94550, USA.

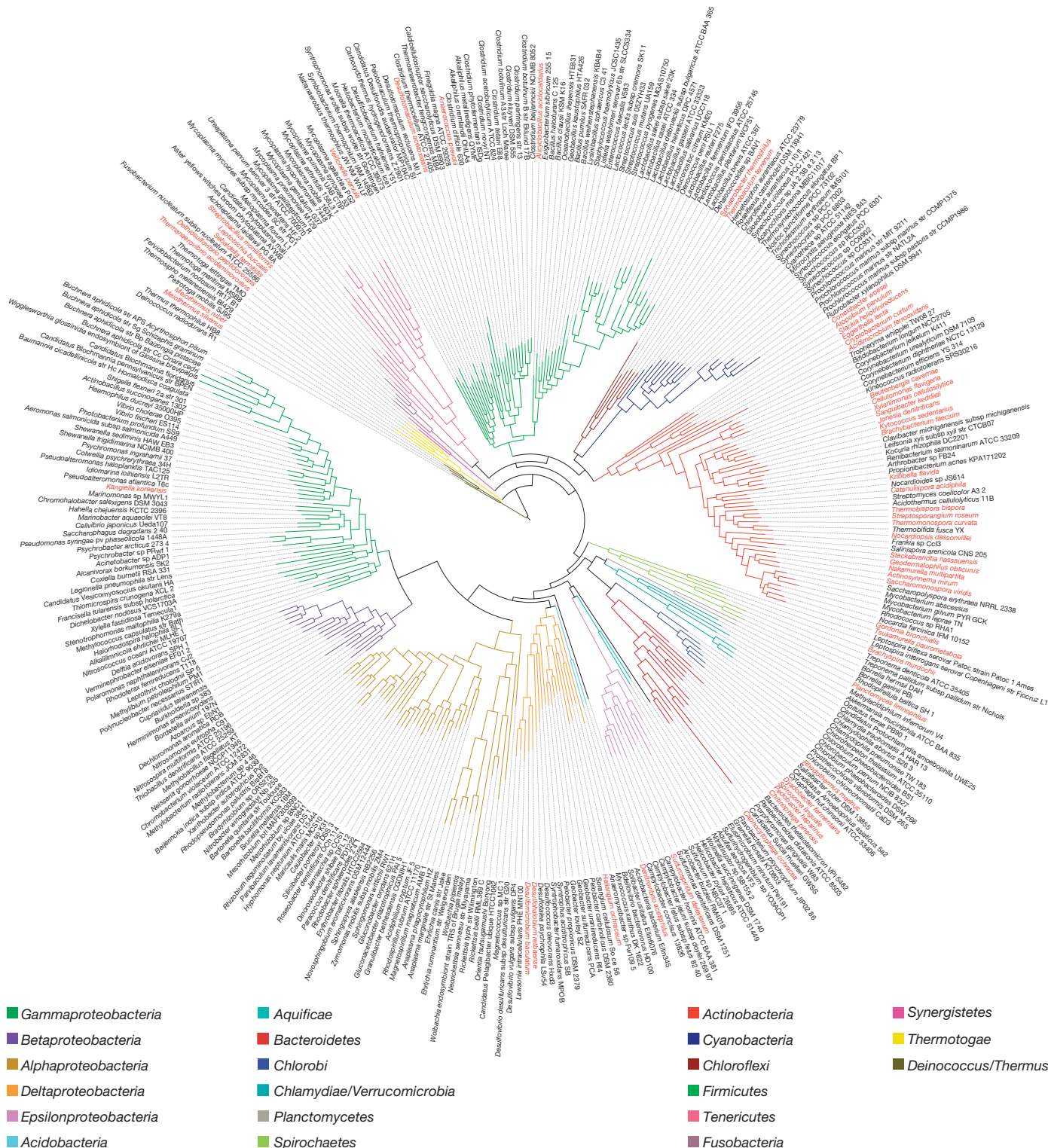


Figure 1 | Maximum-likelihood phylogenetic tree of the bacterial domain based on a concatenated alignment of 31 broadly conserved protein-coding genes¹⁶. Phyla are distinguished by colour of the branch and GEBA genomes are indicated in red in the outer circle of species names.

The discovery and characterization of new gene families and their associated novel functions provide one incentive for sequencing additional genomes, analysis of which has helped to redefine the protein family universe¹⁸. We explored the quantitative effect of tree-based genome selection on the pace of discovery of novel proteins and functions. Specifically, we compared the rate of discovery of novel protein families when progressively adding more closely related genomes versus when adding more distantly related ones (Fig. 2). Granted, many factors contribute to protein family diversity, such

as ecological niche; nevertheless, higher rates of novel protein family discovery were found in the more phylogenetically diverse taxa (Fig. 2). In addition, of the 16,797 families identified in the 56 GEBA genomes, 1,768 showed no significant sequence similarity to any proteins, indicating the presence of novel functional diversity. These results highlight the utility of tree-based genome selection as a means to maximize the identification of novel protein families and argues against lateral gene transfer significantly redistributing genetic novelty between distantly related lineages.

Table 1 | Effect of SSU rRNA tree-based selection of organisms on comparative genomic metrics

Comparative genomic metric	GEBA set	Random sets (number of resamplings)	Fold improvement
Genome tree phylogenetic diversity ¹⁷			
Bacteria (domain)	11.0	3.2 ± 0.7 (100)	2.8–4.4
Actinobacteria (phylum)	4.3	1.4 ± 0.3 (100)	2.5–3.9
New protein family links	46	3 ± 4 (5)	6.6 to >15.3
Genes in new chromosomal cassettes	71,579	16,579 ± 5,523 (20)	3.2–6.5
New gene fusions	433	65 ± 31 (20)	4.5–12.7

GEBA genomes were compared to equivalently sized random sets of reference genomes to quantify the effect of phylogenetic selection.

Novel proteins also can serve to link distantly related homologues whose relatedness would otherwise go undetected. Forty-six such links were identified in the 56 GEBA genomes compared to an average of only three new links in equivalent sets of randomly sampled non-GEBA genomes (Table 1). A useful complement to homology-based predictions of gene function are 'non-homology methods' (ref. 19) such as gene context-based inference that relies on the conserved clustering of functionally related genes across multiple genomes, often in operons or as gene fusions²⁰. We identified over 70,000 genes in new chromosomal cassettes of two or more genes in the GEBA genomes. This represents a three- to sixfold increase over equivalent sets of non-GEBA genomes (Table 1). Similarly, the number of new gene fusions identified in the GEBA genomes is 4 to ~13 times greater than in randomly selected genome sets (Table 1). Because the GEBA data set produced a several-fold improvement over random sets for all metrics examined (Table 1), we predict that other aspects of sequence-based biological discovery will similarly benefit from tree-based genome sequencing.

The GEBA genomes also show significant phylogenetic expansions within known protein families. For example, although only two of the 56 GEBA organisms are known cellulose degraders, we identified in the set of genomes a variety of glycoside hydrolase (GH) genes that may participate in the breakdown of cellulose and hemicelluloses. Among these are 28 and 7 phylogenetically divergent members of the endoglucanase- and processive exoglucanase-containing GH6 and GH48

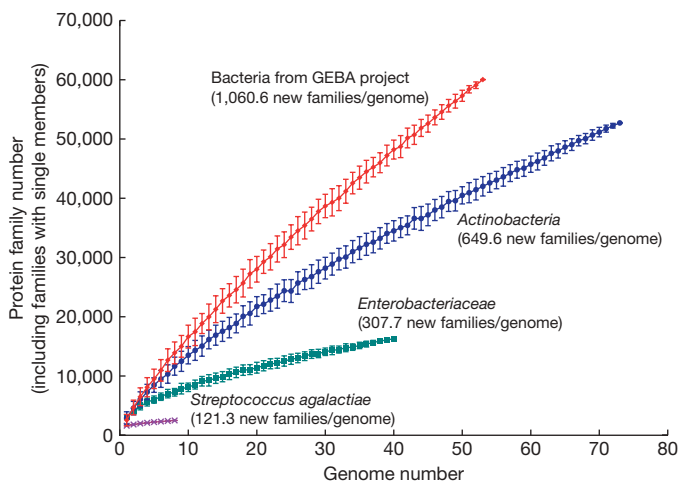


Figure 2 | Rate of discovery of protein families as a function of phylogenetic breadth of genomes. For each of four groupings (species, different strains of *Streptococcus agalactiae*; family, *Enterobacteriaceae*; phylum, *Actinobacteria*; domain, GEBA bacteria), all proteins from that group were compared to each other to identify protein families. Then the total number of protein families was calculated as genomes were progressively sampled from the group (starting with one genome until all were sampled). This was done multiple times for each of the four groups using random starting seeds; the average and standard deviation were then plotted.

1058

families, respectively. *Halorhabdus utahensis*, a halophilic archaeon known to have β -xylanase and β -xylosidase activities²¹, has a chromosomal cluster including two GH10 family β -xylanases and six novel GH5 family proteins of unknown specificity.

The enrichment of genetic diversity is also seen within families of non-coding RNAs, transposable elements, and other cellular components. For example, the genome of the marine myxobacterium *Haliangium ochraceum* contains 807 CRISPR (clustered regularly interspaced short palindromic repeats) units including the largest single CRISPR array known, comprising 382 spacer/repeat units. CRISPR is a newly recognized, but ancient and widespread, system in bacteria and archaea that confers resistance to viruses and other invading foreign DNAs²².

Results from the GEBA pilot project challenge our current understanding for the taxonomic distribution of known gene families. The most striking example of which is the discovery of an actin homologue in *H. ochraceum*. Actin and its close relatives are structural components of the eukaryotic cytoskeleton that are found in every eukaryote and only in eukaryotes. Bacteria and archaea encode instead the shape-determining protein MreB. Although MreBs have some functional and structural similarities to eukaryotic actins, they are regarded, at best, distantly related homologues²³ and possibly not

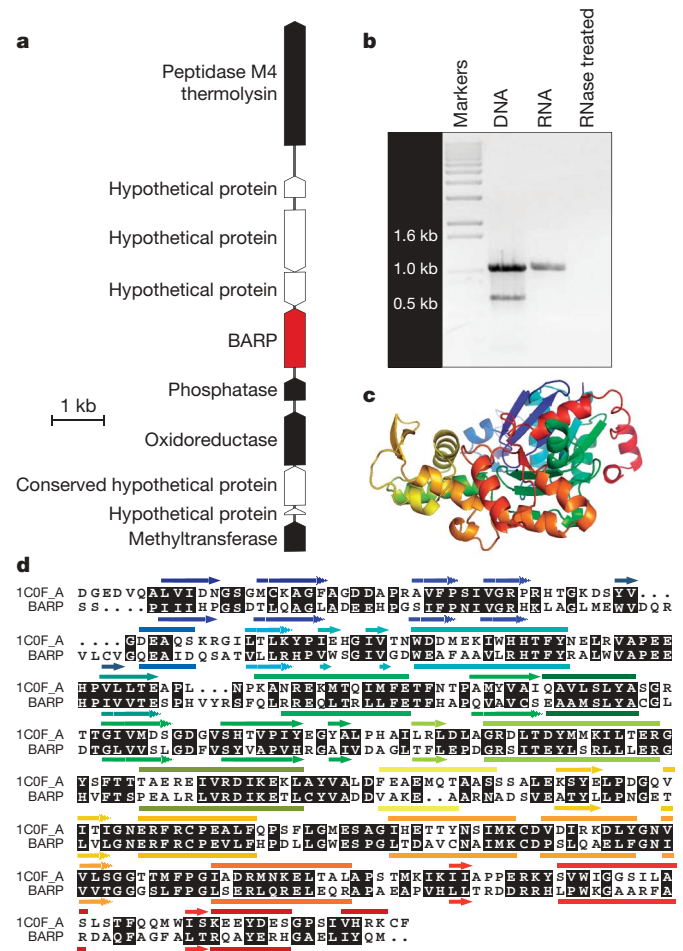


Figure 3 | A bacterial homologue of actin. **a**, Genomic context of the bacterial actin-related protein (BARP) gene within the genome of the marine Deltaproteobacterium *H. ochraceum*. Red, gene encoding BARP; white, genes encoding hypothetical proteins; black, genes with functional annotations. **b**, RT-PCR demonstration of expression of the gene encoding BARP in *H. ochraceum*. **c**, Ribbon plot of the putative structure of BARP. **d**, Alignment of BARP with actin from *Dictyostelium discoideum*²⁹ with similarities in black shaded text. Secondary structure elements (arrows, beta-strands; bars, alpha-helices) are colour-coded as in **c**. A phylogenetic tree including this protein is in Supplementary Figure 1.

even homologous. Like other bacteria, *H. ochraceum* encodes a bona fide MreB protein, but in addition, it encodes a protein that is clearly a member of the actin family, which we have named BARP (bacterial actin-related protein; Fig. 3). Although we do not yet have evidence for its precise function, BARP is expressed in *H. ochraceum* (Fig. 3b). Assuming that the *H. ochraceum mreB* orthologue performs the same function as in other bacteria, and given that the myxobacteria, to which this species belongs, are known to synthesize actin-targeting toxins²⁴, we propose that this BARP may be a dominant-negative inhibitor of eukaryotic actin polymerization. Regardless of its precise function, this first—and so far only—discovery of an expressed homologue of eukaryotic actin in a member of the Bacteria highlights the potential for novel and surprising biological discoveries given a wider genomic sampling of the tree of life.

We conclude that targeting microorganisms for genome sequencing solely on the basis of phylogenetic considerations offers significant far-reaching benefits in diverse areas. Furthermore, the benefits of phylogenetically driven genome sequencing show no sign of saturating with these first 56 genomes. A key question then lies in determining how much bacterial and archaeal diversity remains to be sampled. Using SSU rRNA gene sequences as a proxy for organismal diversity (Fig. 4), we estimate that sequencing the genomes of only 1,520 phylogenetically selected isolates could encompass half of the phylogenetic diversity represented by known cultured bacteria and archaea. Given the continuing reductions in both the cost and difficulty in sequencing genomes²⁵, this is certainly a tractable target in the next few years.

However, the great majority of recognized bacterial and archaeal diversity is not represented by pure cultures and an additional 9,218 genome sequences from currently uncultured species would be required to capture 50% of this recognized diversity (Fig. 4). Such an undertaking will require new approaches to culturing or processing of multi-species samples using methods such as metagenomics²⁶ or physical isolation of cells from mixed populations followed by whole genome amplification methods²⁷. Obtaining reference genomes for the uncultured microbial majority will be a natural extension of the GEBA project, the ultimate goal of which is to provide a phylogenetically balanced genomic representation of the microbial

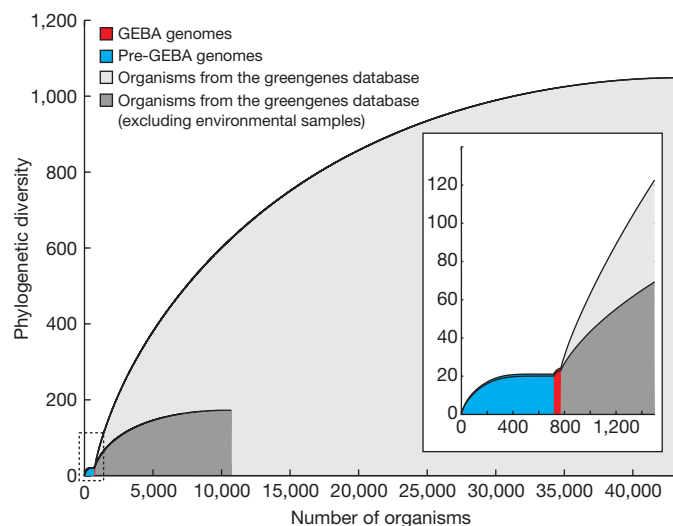


Figure 4 | Phylogenetic diversity of bacteria and archaea on the basis of SSU rRNA genes. Using a phylogenetic tree of unique SSU rRNA gene sequences⁷, phylogenetic diversity was measured for four subsets of this tree: organisms with sequenced genomes pre-GEBA (blue), the GEBA organisms (red), all cultured organisms (dark grey), and all available SSU rRNA genes (light grey). For each subtree, taxa were sorted by their contribution to the subtree phylogenetic diversity³⁰ and the cumulative phylogenetic diversity was plotted from maximal (left) to the least (right). The inset magnifies the first 1,500 organisms. Comparison of the plots shows the phylogenetic ‘dark matter’ left to be sampled.

tree of life. The pilot study presented here is a dedicated first step in this direction.

METHODS SUMMARY

Starting with a phylogenetic tree of SSU rRNA genes⁷, we identified major branches that had no available genome sequences but for which cultured isolates were available in the DSMZ or ATCC culture collections. Selected isolates (Supplementary Table 1a, b) from these branches were grown and DNA isolated (Supplementary Table 1c) and quality checked. DNA was then used for shotgun genome sequencing by Sanger/ABI, Roche/454 and/or Illumina/Solexa technologies (Supplementary Table 2). Sequence reads were assembled separately with different assembly methods and the best draft assembly was used for annotation and as a starting point for genome completion (current genome status is in Supplementary Table 2). Annotation (gene identification, functional prediction, etc.) was performed using the IMG system (<http://img.jgi.doe.gov/gebra>); this was done both after shotgun sequencing and again after genome completion. For ‘whole genome tree’ analysis, a PHYML maximum likelihood phylogenetic tree of a concatenated alignment of 31 marker genes was built using AMPHORA¹⁶. Phylogenetic diversity was calculated as the sum of branch lengths in this and other trees. Protein families were built for various genome sets by using the Markov clustering algorithm (MCL)²⁸ to group proteins on the basis of ‘all versus all’ blastp searches. For analysis of phylogenetic diversity of organisms, a phylogenetic tree was built for a combined alignment of SSU rRNA sequences from published genomes and a non-redundant subset of greengenes SSU rRNA⁷. Further analysis of the genomes was done using IMG database queries and new computational analyses as described in the main text, legends and Supplementary Methods.

Received 3 June; accepted 30 October 2009.

- Fraser, C. M., Eisen, J. A. & Salzberg, S. L. Microbial genome sequencing. *Nature* **406**, 799–803 (2000).
- Liolios, K., Mavromatis, K., Tavernarakis, N. & Kyrpides, N. C. The Genomes On Line Database (GOLD) in 2007: status of genomic and metagenomic projects and their associated metadata. *Nucleic Acids Res.* **36** (database issue), D475–D479 (2008).
- Hugenholtz, P. Exploring prokaryotic diversity in the genomic era. *Genome Biol.* **3**, REVIEWS0003.1–REVIEWS0003.8 (2002).
- Eisen, J. A. Assessing evolutionary relationships among microbes from whole-genome analysis. *Curr. Opin. Microbiol.* **3**, 475–480 (2000).
- Wu, D. *et al.* Complete genome sequence of the aerobic CO-oxidizing thermophile *Thermomicrobium roseum*. *PLoS One* **4**, e4207 (2009).
- Pace, N. R. A molecular view of microbial diversity and the biosphere. *Science* **276**, 734–740 (1997).
- DeSantis, T. Z. *et al.* Greengenes, a chimera-checked 16S rRNA gene database and workbench compatible with ARB. *Appl. Environ. Microbiol.* **72**, 5069–5072 (2006).
- Bernal, A., Ear, U. & Kyrpides, N. Genomes OnLine Database (GOLD): a monitor of genome projects world-wide. *Nucleic Acids Res.* **29**, 126–127 (2001).
- Lapage, S. P. *et al.*, International Code of Nomenclature of Bacteria, 1990 Revision. (American Society for Microbiology, 1992).
- Ward, N., Eisen, J., Fraser, C. & Stackebrandt, E. Sequenced strains must be saved from extinction. *Nature* **414**, 148 (2001).
- Hugenholtz, P. & Kyrpides, N. C. A changing of the guard. *Environ. Microbiol.* **11**, 551–553 (2009).
- Field, D. *et al.* The minimum information about a genome sequence (MIGS) specification. *Nature Biotechnol.* **26**, 541–547 (2008).
- Markowitz, V. M. *et al.* The integrated microbial genomes (IMG) system in 2007: data content and analysis tool extensions. *Nucleic Acids Res.* **36** (database issue), D528–D533 (2008).
- Achtman, M. & Wagner, M. Microbial diversity and the genetic nature of microbial species. *Nature Rev. Microbiol.* **6**, 431–440 (2008).
- Beiko, R. G., Doolittle, W. F. & Charlebois, R. L. The impact of reticulate evolution on genome phylogeny. *Syst. Biol.* **57**, 844–856 (2008).
- Wu, M. & Eisen, J. A. A simple, fast, and accurate method of phylogenomic inference. *Genome Biol.* **9**, R151 (2008).
- Pardi, F. & Goldman, N. Resource-aware taxon selection for maximizing phylogenetic diversity. *Syst. Biol.* **56**, 431–444 (2007).
- Kunin, V., Cases, I., Enright, A. J., de Lorenzo, V. & Ouzounis, C. A. Myriads of protein families, and still counting. *Genome Biol.* **4**, 401 (2003).
- Marcotte, E. M. *et al.* Detecting protein function and protein-protein interactions from genome sequences. *Science* **285**, 751–753 (1999).
- Enright, A. J., Iliopoulos, I., Kyrpides, N. C. & Ouzounis, C. A. Protein interaction maps for complete genomes based on gene fusion events. *Nature* **402**, 86–90 (1999).
- Wainø, M. & Ingvorsen, K. Production of β -xylanase and β -xylosidase by the extremely halophilic archaeon *Haloerhabdus utahensis*. *Extremophiles* **7**, 87–93 (2003).
- Barrangou, R. *et al.* CRISPR provides acquired resistance against viruses in prokaryotes. *Science* **315**, 1709–1712 (2007).
- Doolittle, R. F. & York, A. L. Bacterial actins? An evolutionary perspective. *Bioessays* **24**, 293–296 (2002).

