

Article

Annotation of Protein Domains Reveals Remarkable Conservation in the Functional Make up of Proteomes Across Superkingdoms

Arshan Nasir¹, Aisha Naeem², Muhammad Jawad Khan², Horacio D. Lopez-Nicora³ and Gustavo Caetano-Anollés^{1,*}

¹ Evolutionary Bioinformatics Laboratory, Department of Crop Sciences, University of Illinois, Urbana, IL 61801, USA; E-Mail: anasir@illinois.edu

² Mammalian NutriPhysioGenomics Laboratory, Department of Animal Sciences, University of Illinois, Urbana, IL 61801, USA; E-Mails: naeem1@illinois.edu (A.Na.); khan41@illinois.edu (M.J.K.)

³ Plant Pathology Laboratory, Department of Crop Sciences, University of Illinois, Urbana, IL 61801, USA; E-Mail: hlopezn2@illinois.edu

* Author to whom correspondence should be addressed; E-Mail: gca@illinois.edu; Tel.: +1-217-333-8172; Fax: +1-217-333-8046.

Received: 16 September 2011; in revised form: 28 October 2011 / Accepted: 28 October 2011 /

Published: 8 November 2011

Abstract: The functional repertoire of a cell is largely embodied in its proteome, the collection of proteins encoded in the genome of an organism. The molecular functions of proteins are the direct consequence of their structure and structure can be inferred from sequence using hidden Markov models of structural recognition. Here we analyze the functional annotation of protein domain structures in almost a thousand sequenced genomes, exploring the functional and structural diversity of proteomes. We find there is a remarkable conservation in the distribution of domains with respect to the molecular functions they perform in the three superkingdoms of life. In general, most of the protein repertoire is spent in functions related to metabolic processes but there are significant differences in the usage of domains for regulatory and extra-cellular processes both within and between superkingdoms. Our results support the hypotheses that the proteomes of superkingdom Eukarya evolved via genome expansion mechanisms that were directed towards innovating new domain architectures for regulatory and extra/intracellular process functions needed for example to maintain the integrity of multicellular structure or to

interact with environmental biotic and abiotic factors (e.g., cell signaling and adhesion, immune responses, and toxin production). Proteomes of microbial superkingdoms Archaea and Bacteria retained fewer numbers of domains and maintained simple and smaller protein repertoires. Viruses appear to play an important role in the evolution of superkingdoms. We finally identify few genomic outliers that deviate significantly from the conserved functional design. These include *Nanoarchaeum equitans*, proteobacterial symbionts of insects with extremely reduced genomes, *Tenericutes* and *Guillardia theta*. These organisms spend most of their domains on information functions, including translation and transcription, rather than on metabolism and harbor a domain repertoire characteristic of parasitic organisms. In contrast, the functional repertoire of the proteomes of the Planctomycetes-Verrucomicrobia-Chlamydiae superphylum was no different than the rest of bacteria, failing to support claims of them representing a separate superkingdom. In turn, Protista and Bacteria shared similar functional distribution patterns suggesting an ancestral evolutionary link between these groups.

Keywords: functional annotation; fold superfamily; molecular function; protein domain; SCOP; structure; superkingdom

1. Introduction

Proteins are active components of molecular machinery that perform vital functions for cellular and organismal life [1,2]. Information in the DNA is copied into messenger RNA that is generally translated into proteins by the ribosome. Nascent polypeptide chains are unfolded random coils but quickly undergo conformational changes to produce characteristic and functional folds. These folds are three-dimensional (3D) structures that define the native state of proteins [3,4]. Biologically active proteins are made up of well-packed structural and functional units referred to as domains. Domains appear either singly or in combination with other domains in a protein and act as modules by engaging in combinatorial interplays that enhance the functional repertoires of cells [5]. While molecular interactions between domains in multidomain proteins play important roles in the evolution of protein repertoires [6], it is the domain structure that is maintained in proteins for long periods of evolutionary time [7–9]. This is in sharp contrast to amino acid sequence, which is highly variable. For this reason, protein domains are also considered evolutionary units [7,10–12].

1.1. Classification of Domains

Domains that are evolutionarily related can be grouped together in hierarchical classifications [1,10,13]. One scheme of classifying protein domains is the well-established “Structural Classification of Proteins” (SCOP). The SCOP database groups domains that have sequence conservation (generally with >30% pairwise amino acid residue identities) into fold families (FFs), FFs with structural and functional evidence of common ancestry into fold superfamilies (FSFs), FSFs with common 3D structural topologies into folds (Fs), and Fs sharing a same general architecture into protein classes

[10,14]. SCOP identifies protein domains using concise classification strings (css) (e.g., c.26.1.2, where c represents the protein class, 26 the F, 1 the FSF and 2 the FF). The 97,178 domains indexed in SCOP 1.73 (corresponding to 34,494 PDB entries) are classified into 1,086 F, 1,777 FSFs, and 3,464 FFs. Compared to the number of protein entries in UniProt (531,473 total entries as of July 27, 2011) the number of domain structural designs at these different levels of structural abstraction is quite limited. Their relatively small number suggests that fold space is finite and is evolutionarily highly conserved [1,7,15].

1.2. Assigning FSF Structures to Proteomes

Genome-encoded proteins can be scanned against advanced linear hidden Markov models (HMMs) of structural recognition in SUPERFAMILY [16,17]. HMM libraries are generated using the iterative Sequence Alignment and Modeling (SAM) method. SAM is considered one of the most powerful algorithms for detecting remote homologies [18]. The SUPERFAMILY database currently provides FSF structural assignments for a total of 1,245 model organisms including 96 Archaea, 861 Bacteria and 288 Eukarya.

1.3. Assigning Functional Categories to Protein Domains

Assigning molecular functions to FSFs is a difficult task since approximately 80% of the FSFs defined in SCOP are multi-functional and highly diverse [19]. For example, most of the ancient FSFs, such as the P-loop-containing NTP hydrolase FSF (c.37.1), are highly abundant in nature and include many FFs (20 in case of c.37.1). Each of those families may have functions that impinge on multiple and distinct pathways or networks. The functional annotation scheme introduced by Vogel and Chothia in SUPERFAMILY is a one-to-one mapping scheme that is based on information from various resources, including the Cluster of Orthologous Groups (COG) and Gene Ontology (GO) databases and manual surveys [20–23]. When a FSF is involved in multiple functions, the most predominant function is assigned to that multi-functional FSF under the assumption that the most dominant function is the most ancient and predominantly present in all proteomes. The error rate in assignments is estimated to be <10% for large FSFs and <20% for all FSFs [23].

The SUPERFAMILY functional classification maps seven general functional categories to 50 detailed functional categories in a two-tier hierarchy (Table 1). The seven general categories include *Metabolism*, *Information*, *Intracellular processes (ICP)*, *Extracellular processes (ECP)*, *Regulation*, *General*, and *Other* (we will refer to them as “categories” and “functional repertoires” interchangeably). In this study, we take advantage of this coarse-grained functional annotation scheme to assign individual functional categories to FSFs. We are aware that this one-to-one mapping may not provide a complete profile for multi-functional domains [19]. Dissection of such detailed functions and their comparison across organisms is a difficult problem that we will not address in this study. In contrast, we focus on domains defined at FSF level and use the coarse-grained functional annotation scheme to explore the functional diversity of the proteomes encoded in genomes that have been completely sequenced. Our results yield a global picture of the functional organization of proteomes that is only possible with this classification scheme. Results suggest that the functional structure of proteomes is remarkably conserved across all organisms, ranging from small bacteria to complex

eukaryotes. There is also evidence for the existence of few outliers that deviate from global trends. Here we explore what makes these proteomes distinct.

Table 1. Mapping between the general and minor functional categories for 1,781 protein domains defined in structural classification of proteins (SCOP) 1.73 and the number of fold superfamilies (FSFs) corresponding to each minor category in our dataset of 965 organisms. A total of 135 FSFs could not be annotated. m/tr, metabolism and transport.

Functional category	Minor categories	No. of FSF domains
<i>Metabolism</i> (533 FSFs)	Energy	54
	Photosynthesis	20
	E- transfer	31
	Amino acids m/tr	20
	Nitrogen m/tr	1
	Nucleotide m/tr	30
	Carbohydrate m/tr	30
	Polysaccharide m/tr	21
	Storage	0
	Coenzyme m/tr	50
	Lipid m/tr	17
	Cell envelope m/tr	8
	Secondary metabolism	11
	Redox	55
	Transferases	29
	Other enzymes	156
	<i>General</i> (131 FSFs)	Small molecule binding
Ion binding		13
Lipid/membrane binding		4
Ligand binding		3
General		28
Protein interaction		49
Structural protein		7
<i>Information</i> (201 FSFs)	Chromatin structure	7
	Translation	92
	Transcription	24
	DNA replication/repair	68
	RNA processing	10
	Nuclear structure	0
<i>Other</i> (273 FSFs)	Unknown function	200
	Viral proteins	73
<i>Extracellular processes</i> (95 FSFs)	Cell adhesion	31
	Immune response	19
	Blood clotting	5
	Toxins/defense	40

Table 1. Cont.

Functional category	Minor categories	No. of FSF domains
<i>Intracellular processes</i> (208 FSFs)	Cell cycle, Apoptosis	20
	Phospholipid m/tr	6
	Cell motility	20
	Trafficking/secretion	0
	Protein modification	35
	Proteases	52
	Ion m/tr	21
	Transport	54
<i>Regulation</i> (205 FSFs)	RNA binding, m/tr	19
	DNA-binding	66
	Kinases/phosphatases	15
	Signal transduction	53
	Other regulatory function	34
	Receptor activity	18

2. Results and Discussion

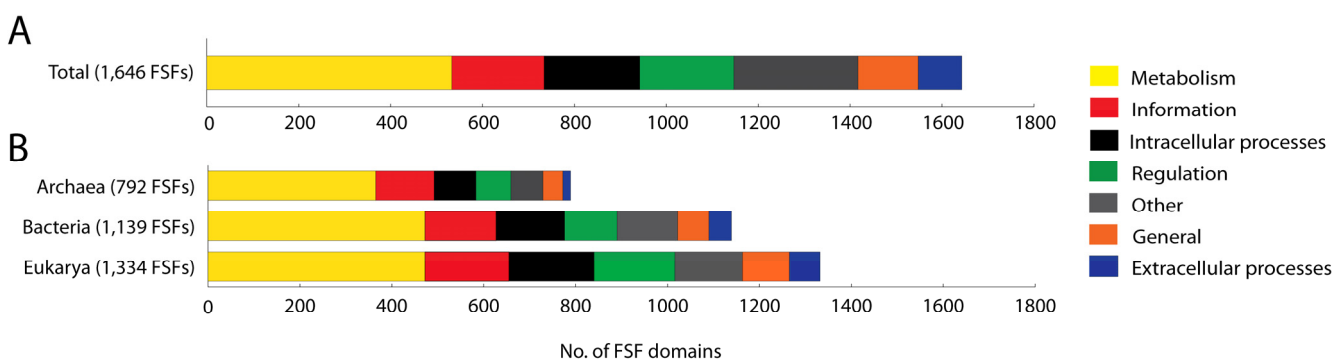
2.1. General Patterns in the Distribution of FSF Domain Functions

We studied the molecular functions of 1,646 domains defined at the FSF level of structural abstraction (SCOP 1.73) that are present in the proteomes of a total of 965 organisms spanning the three superkingdoms. A total of 135 FSFs that could not be annotated were excluded from analysis. For these FSFs, the functional annotation is not available. Out of the 1,646 FSFs studied, approximately one-third (32.38%) performs molecular functions related to *Metabolism*. Categories *Other* (16.58%), *ICP* (12.63%), *Regulation* (12.45%), and *Information* (12.21%) are uniformly distributed within proteomes. In contrast, *General* (7.96%) and *ECP* (5.77%) are significantly underrepresented compared to the rest (Figure 1(A)). The total number of FSFs in each category exhibits the following decreasing trend: *Metabolism* > *Other* > *ICP* > *Regulation* > *Information* > *General* > *ECP*. These patterns of FSF number and relative proteome content are for the most part maintained when studying the functional annotation of FSFs belonging to each superkingdom (Figure 1(B)). However, the number of FSFs in each superkingdom varies considerably and increases in the order Archaea, Bacteria and Eukarya, as we have shown in earlier studies [7].

The significantly higher number of FSFs devoted to *Metabolism* is an anticipated result given the central importance of metabolic networks. However, the much larger number of FSFs corresponding to *Other* is quite unexpected. The 273 FSFs belonging to this category include 200 and 73 FSFs in sub-categories *unknown functions* and *viral proteins*, respectively. The sub-category *unknown function* includes FSFs for which the functions are either unknown or are unclassifiable. Viruses are defined as simple biological entities that are considered to be “gene poor” relatives of cellular organisms [24]. However, the number of domains belonging to *viral proteins* that are present in cellular organisms makes a noteworthy contribution to the total pool of FSFs (4.43%). Thus, viruses have a much more rich and diverse repertoire of domain structures than previously thought and their

association with cellular life has contributed considerable structural diversity to the proteomic make up (A. Nasir, K.M. Kim and G. Caetano-Anollés, ms. in preparation).

Figure 1. Number of protein FSFs annotated for each functional category defined in SCOP 1.73 (A) and in the three superkingdoms (B). The functional distributions show that coarse-grained functions are conserved across cellular proteomes and *Metabolism* is the most dominant functional category. Numbers in parentheses indicate the total number of FSFs annotated in each dataset. The number of FSFs increases in the order Archaea, Bacteria and Eukarya.



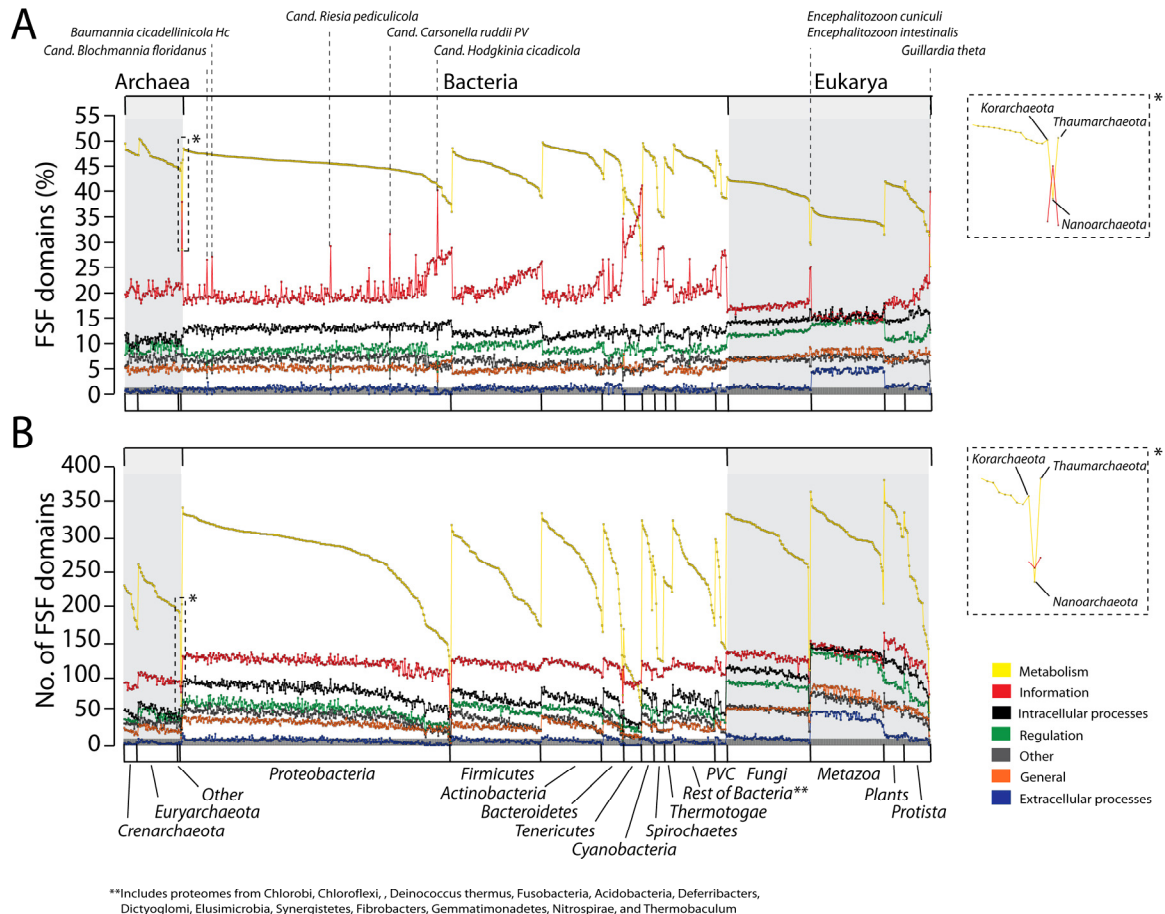
The numbers of FSFs belonging to categories *Regulation*, *Information*, and *ICP* are uniformly distributed in proteomes. However, the *ECP* category is the least represented, perhaps because this category is the last to appear in evolution [7,15]. Extra cellular processes are more important to multicellular organisms (mainly eukaryotes) than to unicellular organisms. Multicellular organisms need efficient communication, such as signaling and cell adhesion. They also trigger immune responses and produce toxins when defending from parasites and pathogens. These *ECP* processes, which are depicted in the minor categories of *cell adhesion*, *immune response*, *blood clotting* and *toxins/defense*, are needed when interacting with environmental biotic and abiotic factors and for maintaining the integrity of multicellular structure. These categories are also present in the microbial superkingdoms but their functional role may be different than in Eukarya.

We note that current genomic research is highly shifted towards the sequencing of microbial genomes, especially those that hold parasitic lifestyles and are of bacterial origin. In fact, 67% of proteomes in our dataset belong to Bacteria. This bias can affect conclusions drawn from global trends such as those in Figure 1(A), including the under-representation of *ECP* FFs, because of their decreased representation in microbial proteomes.

2.2. Distribution of FSF Domain Functions in the Three Superkingdoms of Life

In order to explore whether the overall distribution of general functional categories differs in organisms belonging to the three superkingdoms, we analyzed proteomes at the species level and calculated both the percentage and actual number of FSFs corresponding to different functional repertoires (Figure 2).

Figure 2. The functional distribution of FSFs in individual proteomes of the three superkingdoms. Both the percentage (A) and actual FSF numbers (B) indicate conservation of functional distributions in proteomes and the existence of considerable functional flexibility between superkingdoms. Dotted vertical lines indicate genomic outliers. Insets highlight the interplay between *Metabolism* (yellow trend lines) and *Information* (red trend lines) in *N. equitans*.



FSF domains follow the following decreasing trend in both the percentage and actual counts of FSFs, and do so consistently for the three superkingdoms: *Metabolism* > *Information* > *ICP* > *Regulation* > *Other* > *General* > *ECP*. Note that trend lines across proteomes seldom overlap and cross in Figure 2. It is noteworthy however that this trend differs from the decreasing total numbers of FSFs we described above (Figure 1). Thus, no correlation should be expected between the numbers of FSFs for individual proteomes and the total set for each category. This suggests that variation in functional assignments across proteomes of superkingdoms may not necessarily match overall functional patterns.

Proteomes in microbial superkingdoms Archaea and Bacteria exhibit remarkably similar functional distributions of FSFs (Figure 2(A)). The only exception appears to be the slight overrepresentation of *Regulation* FSFs (green trend lines) and underrepresentation of *ICP* (black trend lines) in Archaea compared to Bacteria (especially Proteobacteria). These distributions are clearly distinct from those in Eukarya. Proteomic representations of FSFs corresponding to *Metabolism* and *Information* are decreased while those of all other five functional categories are significantly and consistently increased

(Figure 2(A)). There is also more variation evident in Eukarya; large groups of proteomes exhibit different patterns of functional use (clearly evident in *Information*; red trend lines in Figure 2(A)).

On the whole, the relative functional make up of the proteomes of individual superkingdoms appear highly conserved (Figure 2(A)). There is however considerable variation in the metabolic functional repertoire of organisms, especially in Bacteria, where *Metabolism* ranges 30–50% of proteomic content (100–350 FSFs, Tables S1 and S2). This variation is not present in other functional repertoires.

Consequently, tendencies of reduction in the metabolic repertoire are generally offset by small increases in the representation of the other six repertoires, with the notable exception of *Information*. In this particular case, when *Metabolism* goes down *Information* goes up. For example, bacterial proteomes with metabolic FSF repertoires of <45% offset their decrease by a corresponding increase in *Information* FSFs (generally from ~20% to ~35%, Figure 2(A)). In all superkingdoms, we identify groups of proteomes or few outliers that deviate from the global trends (vertical dotted lines in Figure 2(A)). As we will discuss below this is generally a consequence of reductive evolution imposed by the lifestyle of organisms (discussed in detail below). Outliers are particularly evident in Bacteria and harbor sharp increases in *Information* repertoires, not always with corresponding decreases in *Metabolism*. In Archaea, decreases of *Metabolism* are generally offset by increases of the *Regulation* category, with an exception in *Nanoarchaeum equitans* (see below). In Eukarya, decreases in *Metabolism* go in hand with decreases in *Information*, and are correspondingly offset mostly by increases in *Regulation* and *ECP*. Apparently, the advantages of regulatory control (e.g., signal transduction and transcriptional and posttranscriptional regulation) and multicellularity counteract the interplay of *Metabolism* and *Information* in eukaryotes.

When we look at the actual number of FSFs within each functional repertoire (Figure 2(B)), we observe a clear trend in domain use that matches the total trend for superkingdoms described above (Figure 1). In most cases, the functional repertoires of Archaea are smaller than those of Bacteria, and bacterial repertoires are generally smaller than those of Eukarya (Figure 2(B)). This holds true for all functional categories. However, the numbers of metabolic FSFs vary 1.5–4 fold in proteomes of superkingdoms, the change being maximal in Bacteria. While both proteomes in Eukarya and Bacteria show similar ranges of metabolic FSFs, the repertoire of Archaea is more constrained. Furthermore, FSFs belonging to categories *Other* and *ECP* are significantly higher in Eukarya than in the microbial superkingdoms. These remarkable observations suggest high conservation in the make up of proteomes of superkingdoms and at the same time considerable levels of flexibility in the metabolic make-up of organisms. Results also support the evolution of the protein complements of Archaea and Bacteria via reductive evolutionary processes and Eukarya by genome expansion mechanisms [7,25]. Reductive tendencies in microbial superkingdoms do not show bias in favor of any functional category. Furthermore, enrichment of eukaryal proteomes with viral proteins supports theories, which state that viruses have played an important role in the evolution of Eukarya [26].

2.3. Distribution of FSF Domain Functions in Individual Phyla/Kingdoms

Figure 2 also describes the functional distribution of FSFs at the phyla/kingdom level for each superkingdom. Plots describing the percentages (Figure 2(A)) and actual number of FSFs in proteomes

(Figure 2(B)) highlight the existence of “outliers” (vertical dotted lines in Figure 2(A)) that deviate from the global functional trends that are typical of each superkingdom.

In Archaea, the functional repertoires of the proteomes of Euryarchaeota, Crenarchaeota, Korarchaeota and Thaumarchaeota were remarkably conserved and consistent with each other. Only *N. equitans* could be considered an outlier (insets of Figure 2). Its proteome deviates from the global archaeal signature by reducing its proteomic make up (it has only 200 distinct FSFs) and by exchanging *Information* for metabolic FSFs. *N. equitans* is an obligate intracellular parasite [27] that is part of a new phylum of Archaea, the Nanoarchaeota [28]. *N. equitans* has many atypical features, including the almost complete absence of operons and presence of split genes [29], tRNA genes that code for only half of the tRNA molecule [30], and the complete absence of the nucleic acid processing enzyme RNase P [31]. Some of these features were used to propose that *N. equitans* is a living fossil [32], represents the root of superkingdom Archaea and the tree of life [33], and is part of a very ancient and yet to be described superkingdom (M. Di Giulio, personal communication). Phylogenomic analyses of domain structures in proteomes suggest Archaea is the most ancient superkingdom [19,34] and has placed *N. equitans* at the base of the tree of life together with other archaeal species. Its ancestral nature is therefore in line with the evolutionary and functional uniqueness of *N. equitans* and the very distinct functional repertoire we here report.

In Bacteria, the functional repertoires of bacterial phyla were also remarkably conserved. Only *Information* and *Metabolism* showed significantly distinct patterns and considerable variation in the use of FSFs. Again, decreases in representation of metabolic FSFs were generally offset by increases in informational FSFs (Figure 2(A)). Notable outliers include the Tenericutes and the Spirochetes. As groups, they have the highest relative usage of *Information* FSFs, which are clearly offset by a decrease in metabolic FSFs. The Tenericutes is a phylum of bacteria that includes class Mollicutes. Members of the Mollicutes are typical obligate parasites of animals and plants (some of medical significance such as *Mycoplasma*) that lack cell walls and have gliding motility. These organisms are characterized by small genome sizes [35] considered to have evolved via reductive evolutionary processes [36]. Because of its unique properties and history, mycoplasmas have been used recently to produce a completely synthetic genome [37]. There were also clear outliers in the Proteobacteria. These included Candidatus *Blochmannia floridanus* (symbiont of ants), *Baumannia cicadellincola* (symbiont of sharpshooter insect), Candidatus *Riesia pediculicola*, Candidatus *Carsonella ruddii* (symbiont of sap-feeding insects) and Candidatus *Hodgkinia cicadicola* (symbiont of cicadas). These bacteria are generally endosymbionts of insects (e.g., ants, sharpshooters, psyllids, cicadas) that have undergone irreversible specialization to an intracellular lifestyle. Candidatus *Carsonella ruddii* has the smallest genome of any bacteria [38]. There were also bacterial proteome groups that were expected to be outliers but were no different than the rest. Bacteria belonging to the superphylum Planctomycetes-Verrucomicrobia-Chlamydiae (PVC) are different from other bacterial phyla because they have an “eukaryotic touch” [39]. Indeed, PVC bacteria display genetic and cellular features that are characteristics of Eukarya and Archaea, including the presence of Histone H1, condensed DNA surrounded by membrane, α -helical repeat domains and β -propeller folds that make up eukaryotic-like membrane coats, reproduction by budding, ether lipids and lack of cell walls [40–42]. Due to the unique nature of the PVC superphylum, it was proposed that these organisms be identified as a separate superkingdom that contributed to the evolution of Eukarya and Archaea [40]. However, trees

of life generated from domain structures in hundreds of proteomes did not dissect the PVC superphylum into a separate group [7,19,34]. Functional distributions of FSFs now show PVC proteomes appear no different from bacteria (Figure 2). These results do not support PVC-inspired theories that explain the diversification of the three cellular superkingdoms of life.

In contrast to the functional repertoires of bacterial and archaeal phyla, proteomes belonging to individual kingdoms in Eukarya had functional signatures that were highly conserved (Figure 2(A)). However, these signatures differed between groups. Plants and fungi had functional representations that were very similar and showed little diversity. In contrast, Metazoa functional distributions increased the representation of *ECP* and *Regulation* FSFs in exchange of FSFs in *Metabolism* and *Information*. Protista had patterns that resemble those of Plants and Fungi but had widely varying metabolic repertoires, very much like Bacteria. This possible link between basal eukaryotes and bacteria revealed by our comparative analysis is consistent with the existence of an ancestor of Bacteria and Eukarya and the early rise of Archaea [34]. Only few outliers belonging to kingdoms Fungi (*Encephalitozoon cuniculi* and *Encephalitozoon intestinalis*) and Protista (*Guillardia theta*) were identified. *E. cuniculi* and *E. intestinalis* are eukaryotic parasites with highly reduced genomes [43,44]. Similarly, *Guillardia theta* is a nucleomorph that has a highly compact and reduced genome with loss of nearly all metabolic genes [45].

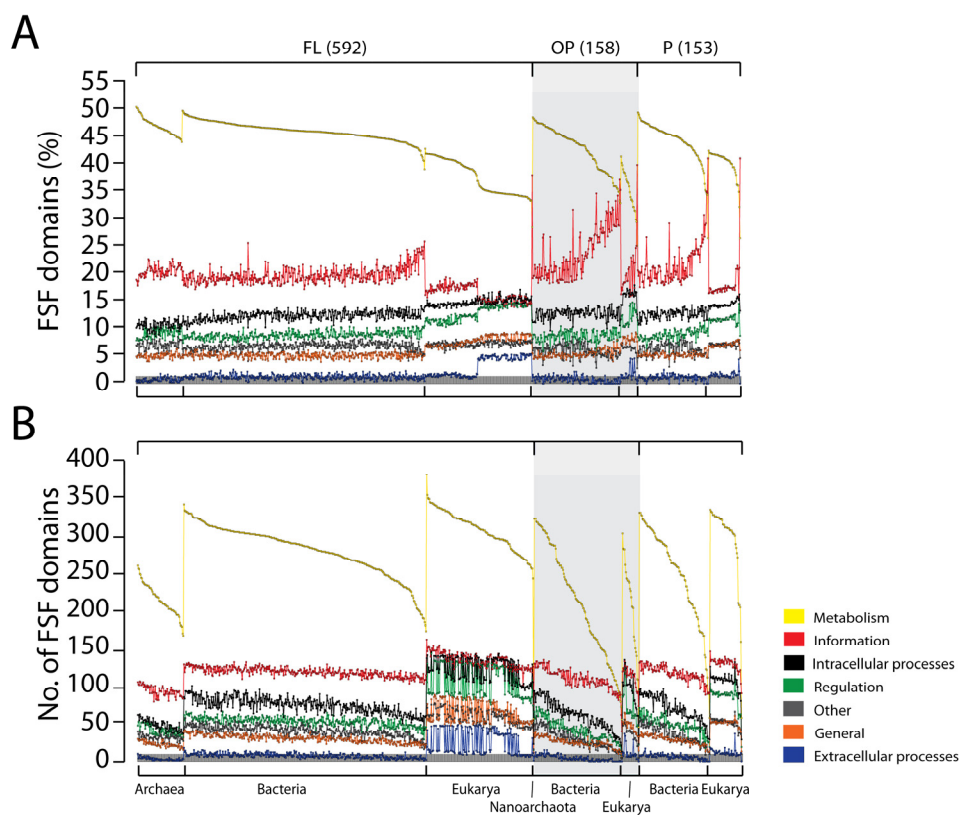
When we look at the actual number of FSFs in proteomes of phyla and kingdoms (Figure 2(B)) we observe that while the overall patterns match those of FSF representation (Figure 2(A)), FSF number revealed considerable variation in the metabolic repertoire of Protista and Bacteria. FSFs in these groups typically ranged 130–340, with PVC and Spirochetes exhibiting the smallest range (130–300 FSFs). In contrast, metabolic repertoires of Archaea and the other eukaryotic kingdoms typically ranged 200–260 FSFs and 270–350 FSFs, respectively. This observation is significant. It provides comparative information to support a unique evolutionary link of phyla within superkingdoms Eukarya and Bacteria. Plots of FSF number also clarified functional patterns in outliers, revealing they did not have more numbers of FSFs in *Information* but rather have reduced metabolic repertoires. This shows that parasitic outliers get rid of metabolic domains and become more and more dependent on host cells.

2.4. Effect of Organism Lifestyle

The analysis thus far revealed the existence of a small group of outliers within each superkingdom. Manual inspection of lifestyles of these organisms showed that all of these organisms are united by a parasitic or symbiotic lifestyle. For example, *N. equitans* is the smallest archaeal genome ever sequenced and represents a new phylum, the Nanoarchaeota [28]. This organism interacts with *Ignicoccus hospitalis*, establishing the only known parasite/symbiont relationship of Archaea, and harbors a highly reduced genome [29]. Parasitic/symbiotic relationships with various plants and animals can be found in *Tenericutes* and in the endosymbionts of insects that belong to *Proteobacteria*. Similarly, the *Encephalitozoon* species are eukaryotic parasites that lack mitochondria and have highly reduced genomes [43,44]. *E. cuniculi* has even a chromosomal dispersion of its ribosomal genes, very much like *N. equitans*, and the rRNA of the large ribosomal subunit reduced to its universal core [46]. Similarly, *Guillardia theta* is a nucleomorph that has a highly compact and reduced genome with loss of nearly all metabolic genes [45]. Thus, all outliers exhibit extreme or unique cases of genome reduction.

In order to explore whether organisms that engage in parasitic or symbiotic interactions have general tendencies that resemble those of the outliers, we classified organisms into three different lifestyles: free living (FL) (592 proteomes), facultative parasitic (P) (153 proteomes), and obligate parasitic (OP) (158 proteomes). Functional distributions for the seven general functional categories for these proteomic sets explained the role of parasitic life on proteomic constitution (Figure 3). Plots of percentages (Figure 3(A)) and actual number of FSFs in proteomes (Figure 3(B)) showed FSF distribution in FL organisms were remarkably homogenous and that the vast majority of variability within superkingdoms was ascribed to the P and OP lifestyles. This variability was for the most part explained by a sharp decline in the number of metabolic FSFs that are assigned to the *Metabolism* general category (Figure 3(B)). Plots also support the hypothesis that parasitic organisms have gone the route of massive genome reduction in a tendency to loose all of their metabolic genes. This tendency makes them more and more dependent on host cells for metabolic functions and survival [47,48].

Figure 3. The functional distribution of FSFs with respect to organism lifestyle. Both the percentage (A) and actual FSF numbers (B) indicate that obligate parasitic (OP) and facultative parasitic (P) organisms exhibit considerable variability in their metabolic repertoires (yellow trend lines) that is offset by corresponding increases in the *Information* FSFs (red trend lines).



The number of domains corresponding to each general functional category in the proteomes of FL organisms increases in the order Archaea, Bacteria and Eukarya (Table S3). When compared to the total proteomic set (Figure 2), *Metabolism* remains the predominant functional category and a large number of domains in all the proteomes perform metabolic functions. Again, the proteomes of Eukarya

have the richest FSF repertoires, and those of Archaea the most simple. Since maximum variability lies within the proteome repertoires of parasitic/symbiotic organisms (Figure 3) and parasitism/symbiosis in these organisms is the result of secondary adaptations, the analysis of proteomic diversity in FL organisms allows us to test if the functional repertoires of superkingdoms are indeed statistically significant. Analysis of variance showed that the number of FSFs for each functional repertoire was consistently different between superkingdoms ($p < 0.0001$; Table S3). This supports the conclusions drawn from earlier analyses that the microbial superkingdoms followed a genome reduction path while Eukarya expanded their genomic repertoires [7,25].

2.5. Analysis of Minor Functional Categories

The seven general categories of molecular functions map to 50 minor categories (Table 1). We explored the distribution of FSFs corresponding to each minor category in superkingdoms (Figure 4). Only category “not annotated” (NONA) was excluded from analysis. In terms of percentage (Figure 4(A)), the overall functional signature is split into two components: prokaryotic and eukaryotic. Prokaryotes spend most of their domain repertoire on *Metabolism* and *Information* whereas Eukarya stand out in *ECP* (particularly cell adhesion, immune response), *Regulation* (DNA binding, signal transduction), and all the minor functional categories corresponding to *ICP* and *General*.

Figure 4. The percentage (A) and number (B) of FSFs in minor functional categories across superkingdoms. Archaea (A) and Bacteria (B) spend most of their proteomes in functions related to *Metabolism* and *Information* whereas Eukarya (E) stand out in the minor categories of *Regulation*, *General*, *Intracellular processes (ICP)* and *Extracellular processes (ECP)*. In turn, the number of FSFs increases in the order Archaea, Bacteria and Eukarya. Eukaryal proteomes have the richest functional repertoires for *Regulation*, *Other*, *General*, *ICP* and *ECP*.

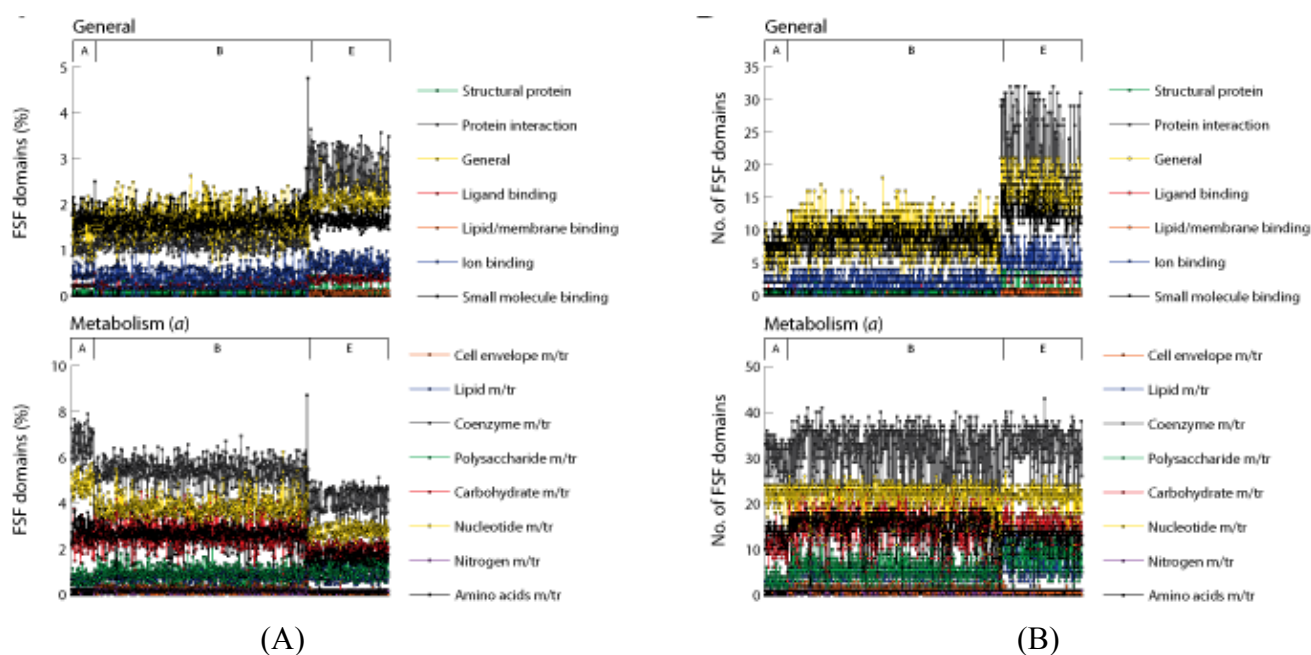
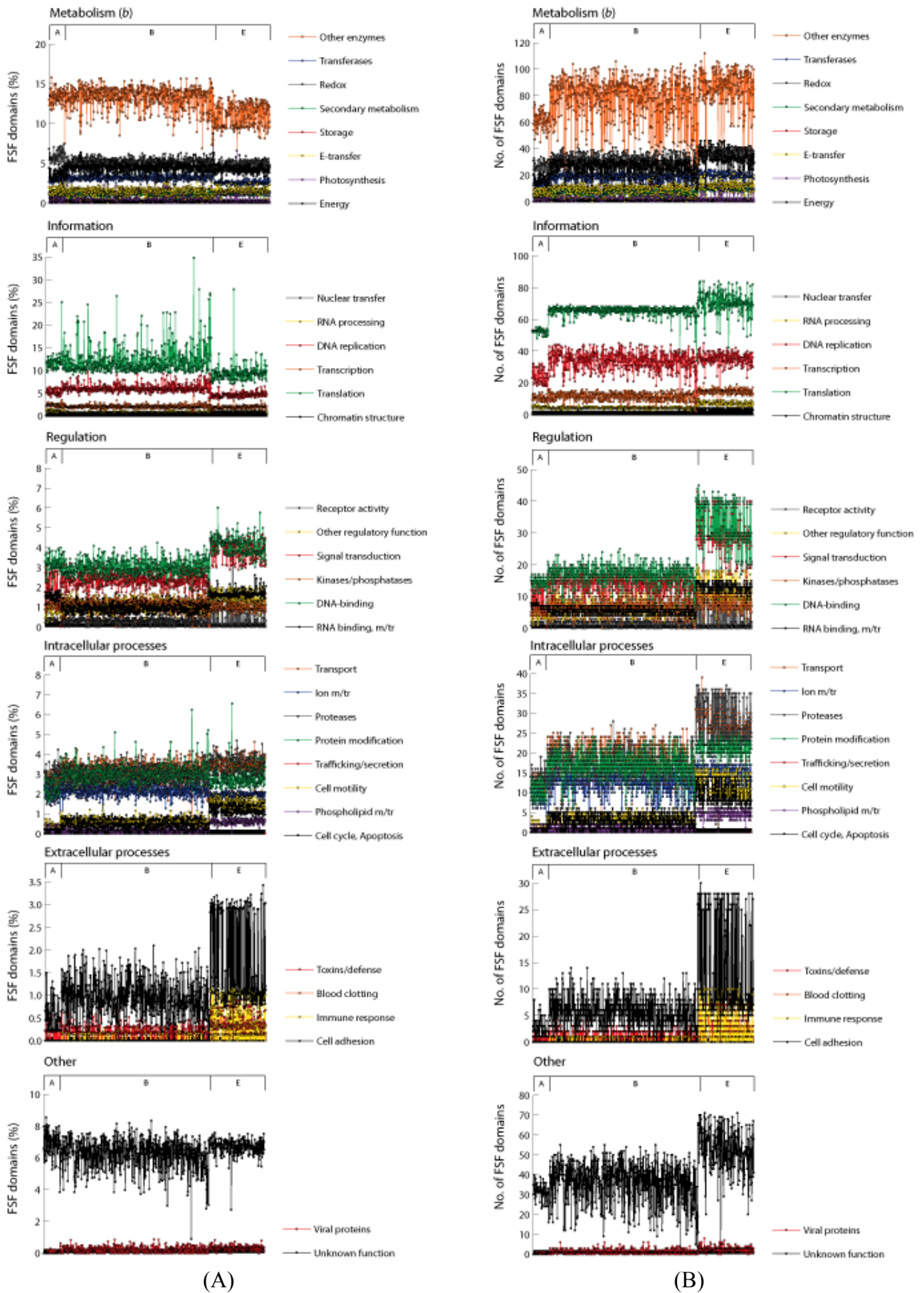


Figure 4. Cont.



In terms of domain counts (Figure 4(B)), proteomes of Eukarya have the richest functional repertoires with a significantly large number of FSFs devoted for each minor functional category. Bacteria and Archaea work with small number of domains. However, the number of FSFs in Bacteria is significantly higher compared to Archaea (supporting results of Figures 1 and 2 and Table S3). These results are consistent with the evolutionary trends in proteomes described previously [7,19,25]. Our results support the complex nature of the Last Universal Common Ancestor (LUCA) [19] and are consistent with the evolution of microbial superkingdoms via reductive evolutionary processes and the evolution of eukaryal proteomes by genome expansion [7,25]. It appears that Archaea went on the route of genome reduction very early in evolution and was followed by Bacteria and finally Eukarya. Late in evolution, the eukaryal superkingdom increased the representation of FSFs and developed a rich proteome. This can explain the relatively huge and diverse nature of eukaryal proteomes compared to prokaryotic proteomes. Finally, there appears to be no significant difference in the distributions of FSFs corresponding to *Metabolism* and *Information* between Bacteria and Eukarya except for minor category “Translation” (green trend lines in Figures 4(B, *Information*)) that is significantly higher in Eukarya compared to Bacteria. This shows that Bacteria exhibit incredible metabolic and informational diversity despite their reduced genomic complements. We conclude that the genome expansion in Eukarya occurred primarily for functions related to *ECP*, *ICP*, *Regulation* and *General*.

2.6. Reliability of Functional Annotations and Conclusions of this Study

Our analysis depends upon the accuracy of assigning structures to protein sequences and the SCOP protein classification and SUPERFAMILY functional annotation schemes. Databases such as SCOP and SUPERFAMILY are continuously updated with more and more genomes and new assignments. We therefore ask the reader to focus on the general trends in the data as opposed to the specifics such as the exact percentage or numbers of FSFs in each functional repertoire. Trends related to the number of domains in Archaea relative to Bacteria and Eukarya and the reduction of metabolic repertoires in parasitic organisms should be considered robust since these have been reliably observed in previous studies with more limited datasets [1,7,15,19,34]. Biases in sampling of proteomes in the three superkingdoms is not expected to over or underestimate the remarkably conserved nature of the functional makeup. We show that the conservation of molecular functions in proteomes is only broken in genomic outliers that are united by parasitic lifestyles. Thus equal sampling will not significantly alter the global trends described for individual superkingdoms. In light of our results, organism lifestyle is the only factor affecting the conserved nature of proteomes. Finally, we propose that lower or higher than expected numbers of FSFs in any category (subcategory) can be explained either by possible limitations of the scheme used to annotate molecular functions of FSFs or the simple nature of the functional repertoire. For example, the number of FSFs in subcategory structural proteins (main category *General*) is 7 (Table 1) despite the importance of structural proteins in cellular organization. Table S4 lists the description of these FSFs and shows that indeed these FSF domains play important structural roles. Their limited number indicates that the structural and functional organization is quite limited and very few folds play important structural roles. Another possibility is the “hidden” overlap between FSFs and molecular functions due to the one-to-one mapping limitations of the SUPERFAMILY functional annotation scheme. Most of the large FSFs include many FFs and

participate in multiple pathways; for few FSFs a complete functional profile may not be intuitively obvious. This may be one of the shortcomings of using this functional annotation scheme but dissection of such detailed functions and pathways is a difficult task and is not described in this study. In summary, we do not believe that the classification or annotation schemes, despite their limitations, would undergo serious revisions or weaken our findings.

3. Experimental Section

3.1. Data Retrieval

We downloaded the protein architecture assignments for a total of 965 organisms including 70 Archaea, 651 Bacteria and 244 Eukarya (Table S5) from SUPERFAMILY ver. 1.73 MySQL [16,17] at an *E*-value cutoff of 10^{-4} . This cutoff is considered a stringent threshold to eliminate the rate of false positives in HMM assignments [19]. Classification of organisms according to their lifestyles was done manually and resulted in 592 FL, 153 P, and 158 OP organisms.

3.2. Assigning Functional Categories to Protein Domains

The most recent domain functional annotation file for SCOP 1.73 was downloaded from the SUPERFAMILY webserver [23]. For each genome we extracted the set of unique FSFs present and then mapped them to the 7 general and 50 detailed functional categories. We calculated both the percentage and actual number of domains using programming implementations in Python 3.1 (<http://www.python.org/download/>).

3.3. Statistical Analysis

The statistical significance between the numbers of functional FSFs in FL organisms of superkingdoms was evaluated by Welch's ANOVA in SAS (<http://www.sas.com/software/sas9>), which is the appropriate test to detect differences between means for groups having unequal variances [49]. We excluded organisms with P and OP lifestyles in order to remove noise from the data. Additionally, in order to meet asymptotic normality, we used the Log_{10} transformation and rescaled the data to 0–7 using the following formula,

$$N_{normal} = [\text{Log}_{10}(N_{xy})/\text{Log}_{10}(N_{max})] \times 7$$

where N_{xy} is the count of a FSF in x functional category in y superkingdom; N_{max} is the largest value in the matrix and N_{normal} is the normalized and scaled score for FSF x in y superkingdom.

4. Conclusions

Our analysis revealed a remarkable conservation in the functional distribution of protein domains in superkingdoms for proteomes for which we have structural assignments. Figure S1 showcases average distribution of FSFs in phyla, kingdoms, and superkingdoms. The biggest proportion of each proteome is devoted in all cases to functions related to *Metabolism*. Phylogenomic analysis has shown that *Metabolism* appeared earlier than other functional groups and their structures were the first to spread in life [1,50]. This would explain the relative large representation of *Metabolism* in the functional toolkit

of cells. Usage of domains related to *ECP* and *Regulation* is significantly higher in *Metazoa* compared to the rest. This showcases the importance of regulation signal transduction mechanisms for eukaryotic organisms [51,52]. Our results support the view that prokaryotes evolved via reductive evolutionary processes whereas genome expansion was the route taken by eukaryotic organisms. Genome expansion in Eukarya seems to be directed towards innovation of FSF architectures, especially those linked to *Regulation*, *ECP* and *General*. Finally, viral structures make up a substantial proportion of cellular proteomes and appear to have played an important role in the evolution of cellular life.

Organisms with parasitic lifestyles have simple and reduced proteomes and rely on host cells for metabolic functions. Tenericutes are unique in this regard. They spend most of their proteomic resources in functions linked to *Information* (e.g., translation, replication). Remarkably, we find that the conservation of molecular functions in proteomes is only broken in “outliers” with parasitic lifestyles that do not obey the global trends. We conclude that organism lifestyle is a crucial factor in shaping the nature of proteomes.

Acknowledgements

This study began as a class project in CPSC 567, a course in bioinformatics and systems biology taught by G.C.-A. at the University of Illinois in spring 2011. We thank Kyung Mo Kim and Liudmila Yafremava for information about lifestyles. A.N., A.Na., M.J.K. and H.D.L.-N. conceived the experiments and analyzed the data. G.C.-A. supervised the project and edited the manuscript. Research was supported by the National Science Foundation (MCB-0749836), CREES-USDA and the Soybean Disease Biotechnology Center (to G.C.-A.). Any opinions, findings, and conclusions and recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the funding agencies.

References

1. Caetano-Anolles, D.; Kim, K.M.; Mittenthal, J.E.; Caetano-Anolles, G. Proteome evolution and the metabolic origins of translation and cellular life. *J. Mol. Evol.* **2011**, *72*, 14–33.
2. Lesk, A.M. *Introduction to Protein Architecture*; Oxford University Press: New York, NY, USA, 2001.
3. Cordes, M.H.; Davidson, A.R.; Sauer, R.T. Sequence space, folding and protein design. *Curr. Opin. Struct. Biol.* **1996**, *6*, 3–10.
4. Linderstrom-Lang, K.U.; Schellman, J.A. *The Enzymes*; Academic Press: New York, NY, USA, 1959; pp. 443–510.
5. Wang, M.; Caetano-Anolles, G. The evolutionary mechanics of domain organization in proteomes and the rise of modularity in the protein world. *Structure* **2009**, *17*, 66–78.
6. Vogel, C.; Bashton, M.; Kerrison, N.D.; Chothia, C.; Teichmann, S.A. Structure, function and evolution of multidomain proteins. *Curr. Opin. Struct. Biol.* **2004**, *14*, 208–216.
7. Wang, M.; Yafremava, L.S.; Caetano-Anolles, D.; Mittenthal, J.E.; Caetano-Anolles, G. Reductive evolution of architectural repertoires in proteomes and the birth of the tripartite world. *Genome Res.* **2007**, *17*, 1572–1585.

8. Gerstein, M.; Hegyi, H. Comparing genomes in terms of protein structure: Surveys of a finite parts list. *FEMS Microbiol. Rev.* **1998**, *22*, 277–304.
9. Chothia, C.; Gough, J.; Vogel, C.; Teichmann, S.A. Evolution of the protein repertoire. *Science* **2003**, *300*, 1701–1703.
10. Murzin, A.G.; Brenner, S.E.; Hubbard, T.; Chothia, C. Scop: A structural classification of proteins database for the investigation of sequences and structures. *J. Mol. Biol.* **1995**, *247*, 536–540.
11. Orengo, C.A.; Michie, A.D.; Jones, S.; Jones, D.T.; Swindells, M.B.; Thornton, J.M. Cath—A hierarchic classification of protein domain structures. *Structure* **1997**, *5*, 1093–1108.
12. Riley, M.; Labedan, B. Protein evolution viewed through escherichia coli protein sequences: Introducing the notion of a structural segment of homology, the module. *J. Mol. Biol.* **1997**, *268*, 857–868.
13. Ponting, C.P.; Russell, R.R. The natural history of protein domains. *Annu. Rev. Biophys. Biomol. Struct.* **2002**, *31*, 45–71.
14. Andreeva, A.; Howorth, D.; Chandonia, J.M.; Brenner, S.E.; Hubbard, T.J.; Chothia, C.; Murzin, A.G. Data growth and its impact on the scop database: New developments. *Nucleic Acids Res.* **2008**, *36*, D419–D425.
15. Caetano-Anolles, G.; Wang, M.; Caetano-Anolles, D.; Mittenthal, J.E. The origin, evolution and structure of the protein world. *Biochem. J.* **2009**, *417*, 621–637.
16. Gough, J.; Karplus, K.; Hughey, R.; Chothia, C. Assignment of homology to genome sequences using a library of hidden markov models that represent all proteins of known structure. *J. Mol. Biol.* **2001**, *313*, 903–919.
17. Wilson, D.; Madera, M.; Vogel, C.; Chothia, C.; Gough, J. The superfamily database in 2007: Families and functions. *Nucleic Acids Res.* **2007**, *35*, D308–D313.
18. Karplus, K. Sam-t08, hmm-based protein structure prediction. *Nucleic Acids Res.* **2009**, *37*, W492–W497.
19. Kim, K.M.; Caetano-Anolles, G. The proteomic complexity and rise of the primordial ancestor of diversified life. *BMC Evol. Biol.* **2011**, *11*, 140:1–140:24.
20. Vogel, C.; Berzuini, C.; Bashton, M.; Gough, J.; Teichmann, S.A. Supra-domains: Evolutionary units larger than single protein domains. *J. Mol. Biol.* **2004**, *336*, 809–823.
21. Vogel, C.; Teichmann, S.A.; Pereira-Leal, J. The relationship between domain duplication and recombination. *J. Mol. Biol.* **2005**, *346*, 355–365.
22. Vogel, C.; Chothia, C. Protein family expansions and biological complexity. *PLoS Comput. Biol.* **2006**, *2*, e48:0370–e48:0382.
23. Vogel, C. Function annotation of SCOP domain superfamilies 1.73. Superfamily-HMM library and genome assignments server. Available online: <http://supfam.cs.bris.ac.uk/SUPERFAMILY/function.html> (accessed on 28 October 2011).
24. Moreira, D.; Lopez-Garcia, P. Ten reasons to exclude viruses from the tree of life. *Nat. Rev. Microbiol.* **2009**, *7*, 306–311.
25. Wang, M.; Kurland, C.G.; Caetano-Anolles, G. Reductive evolution of proteomes and protein structures. *Proc. Natl. Acad. Sci. USA* **2011**, *108*, 11954–11958.
26. Koonin, E.V.; Wolf, Y.I.; Nagasaki, K.; Dolja, V.V. The big bang of picorna-like virus evolution antedates the radiation of eukaryotic supergroups. *Nat. Rev. Microbiol.* **2008**, *6*, 925–939.

27. Das, S.; Paul, S.; Bag, S.K.; Dutta, C. Analysis of nanoarchaeum equitans genome and proteome composition: Indications for hyperthermophilic and parasitic adaptation. *BMC Genomics* **2006**, *7*, 186:1–186:16.
28. Huber, H.; Hohn, M.J.; Rachel, R.; Fuchs, T.; Wimmer, V.C.; Stetter, K.O. A new phylum of archaea represented by a nanosized hyperthermophilic symbiont. *Nature* **2002**, *417*, 63–67.
29. Waters, E.; Hohn, M.J.; Ahel, I.; Graham, D.E.; Adams, M.D.; Barnstead, M.; Beeson, K.Y.; Bibbs, L.; Bolanos, R.; Keller, M.; Kretz, K.; Lin, X.; Mathur, E.; Ni, J.; Podar, M.; Richardson, T.; Sutton, G.G.; Simon, M.; Soll, D.; Stetter, K.O.; Short, J.M.; Noordewier, M. The genome of Nanoarchaeum equitans: Insights into early archaeal evolution and derived parasitism. *Proc. Natl. Acad. Sci. USA* **2003**, *100*, 12984–12988.
30. Randau, L.; Munch, R.; Hohn, M.J.; Jahn, D.; Soll, D. Nanoarchaeum equitans creates functional trnas from separate genes for their 5'- and 3'-halves. *Nature* **2005**, *433*, 537–541.
31. Randau, L.; Schroder, I.; Soll, D. Life without rnase p. *Nature* **2008**, *453*, 120–123.
32. Di Giulio, M. Nanoarchaeum equitans is a living fossil. *J. Theor. Biol.* **2006**, *242*, 257–260.
33. Di Giulio, M. The tree of life might be rooted in the branch leading to nanoarchaeota. *Gene* **2007**, *401*, 108–113.
34. Kim, K.M.; Caetano-Anolles, G. The evolutionary history of protein fold families and proteomes confirms Archaea is the most ancient superkingdom. Ms. submitted.
35. Woese, C.R.; Maniloff, J.; Zablen, L.B. Phylogenetic analysis of the mycoplasmas. *Proc. Natl. Acad. Sci. USA* **1980**, *77*, 494–498.
36. Chambaud, I.; Heilig, R.; Ferris, S.; Barbe, V.; Samson, D.; Galisson, F.; Moszer, I.; Dybvig, K.; Wróblewski, H.; Viari, A.; Rocha, E.P.; Blanchard, A. The complete genome sequence of the murine respiratory pathogen Mycoplasma pulmonis. *Nucleic Acids Res.* **2001**, *29*, 2145–2153.
37. Gibson, D.G.; Smith, H.O.; Hutchison, C.A., III.; Venter, J.C.; Merryman, C. Chemical synthesis of the mouse mitochondrial genome. *Nat. Methods* **2010**, *7*, 901–903.
38. Nakabachi, A.; Yamashita, A.; Toh, H.; Ishikawa, H.; Dunbar, H.E.; Moran, N.A.; Hattori, M. The 160-kilobase genome of the bacterial endosymbiont carsonella. *Science* **2006**, *314*, 267.
39. Forterre, P.; Gribaldo, S. Bacteria with a eukaryotic touch: A glimpse of ancient evolution? *Proc. Natl. Acad. Sci. USA* **2010**, *107*, 12739–12740.
40. Santarella-Mellwig, R.; Franke, J.; Jaedicke, A.; Gorjanacz, M.; Bauer, U.; Budd, A.; Mattaj, I.W.; Devos, D.P. The compartmentalized bacteria of the planctomycetes-verrucomicrobia-chlamydiae superphylum have membrane coat-like proteins. *PLoS Biol.* **2010**, *8*, e1000281:1–e1000281:11.
41. Kamneva, O.K.; Liberles, D.A.; Ward, N.L. Genome-wide influence of indel substitutions on evolution of bacteria of the PVC superphylum, revealed using a novel computational method. *Genome Biol. Evol.* **2010**, *2*, 870–886.
42. Devos, D.P.; Reynaud, E.G. Evolution. Intermediate steps. *Science* **2010**, *330*, 1187–1188.
43. Katinka, M.D.; Duprat, S.; Cornillot, E.; Méténier, G.; Thomarat, F.; Prensier, G.; Barbe, V.; Peyretailade, E.; Brottier, P.; Wincker, P.; Delbac, F.; El Alaoui, H.; Peyret, P.; Saurin, W.; Gouy, M.; Weissenbach, J.; Vivarès, C.P. Genome sequence and gene compaction of the eukaryote parasite Encephalitozoon cuniculi. *Nature* **2001**, *414*, 450–453.

44. Corradi, N.; Pombert, J.F.; Farinelli, L.; Didier, E.S.; Keeling, P.J. The complete sequence of the smallest known nuclear genome from the microsporidian *Encephalitozoon intestinalis*. *Nat. Commun.* **2010**, *1*, 77, doi: 10.1038/ncomms1082.
45. Douglas, S.; Zauner, S.; Fraunholz, M.; Beaton, M.; Penny, S.; Deng, L.T.; Wu, X.; Reith, M.; Cavalier-Smith, T.; Maier, U.G. The highly reduced genome of an enslaved algal nucleus. *Nature* **2001**, *410*, 1091–1096.
46. Peyretailade, E.; Biderre, C.; Peyret, P.; Duffieux, F.; Metenier, G.; Gouy, M.; Michot, B.; Vivares, C.P. Microsporidian encephalitozoon cuniculi, a unicellular eukaryote with an unusual chromosomal dispersion of ribosomal genes and a *lsu rRNA* reduced to the universal core. *Nucleic Acids Res.* **1998**, *26*, 3513–3520.
47. Martin, W.; Herrmann, R.G. Gene transfer from organelles to the nucleus: How much, what happens, and why? *Plant Physiol.* **1998**, *118*, 9–17.
48. Keeling, P.J.; Slamovits, C.H. Causes and effects of nuclear genome reduction. *Curr. Opin. Genet. Dev.* **2005**, *15*, 601–608.
49. Welch, B.L.; The significance of the difference between two means when the population variances are unequal. *Biometrika* **1938**, *29*, 350–362.
50. Caetano-Anolles, G.; Kim, H.S.; Mittenthal, J.E. The origin of modern metabolic networks inferred from phylogenomic analysis of protein architecture. *Proc. Natl. Acad. Sci. USA* **2007**, *104*, 9358–9363.
51. Ingham, P.W.; Nokano, Y.; Seger, C. Mechanisms and functions of Hedgehog signalling across the metazoa. *Nat. Rev. Genet.* **2011**, *12*, 393–406.
52. Bürglin, T.R. Evolution of hedgehog and hedgehog-related genes, their origin from Hog proteins in ancestral eukaryotes and discovery of a novel Hint motif. *BMC Genomics* **2008**, *9*, 127:1–127:28.

Supplementary Materials

Figure S1. Average distribution of FSFs in phyla, kingdom, and superkingdoms suggest conservation of functional design in proteomes. Numbers in parentheses indicate total number of proteomes analyzed for each phyla/kingdom.

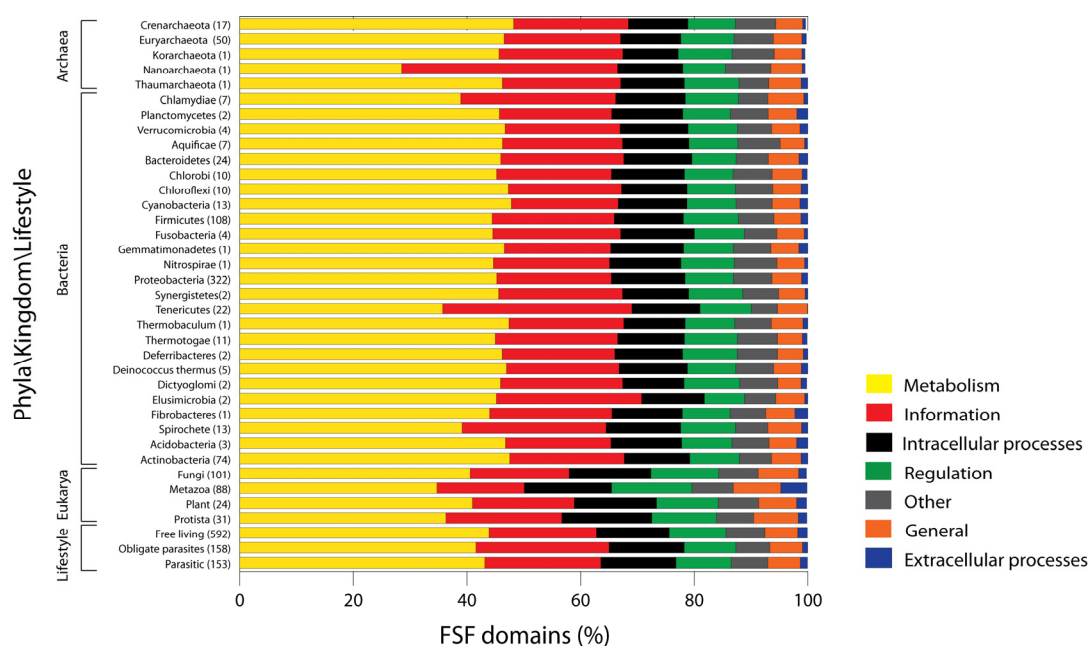


Table S1. Average number of FSF domains in each phyla/kingdom corresponding to the seven general functional categories. Numbers were rounded up when the decimal value exceeded 0.5 and rounded down otherwise. Nanoarchaeota and Tenericutes have the least number of metabolic domains and are highlighted in bold. Eukaryal kingdoms (Fungi, Metazoa, Plants and Protista) have the richest FSF repertoires compared to the prokaryotes.

Superkingdom	Phyla/Kingdom	Metabolism	Information	ICP	Regulation	Other	General	ECP
Archaea	Crenarchaeota	204	85	44	35	30	20	2
	Euryarchaeota	219	96	50	44	32	24	4
	Korarchaeota	178	85	38	37	29	19	2
	Nanoarchaeota	57	76	23	15	16	11	1
	Thaumarchaeota	202	91	49	42	23	25	5
Bacteria	Proteobacteria	274	119	78	52	42	31	7
	Firmicutes	246	117	67	53	35	26	7
	Actinobacteria	275	115	66	50	33	30	7
	Bacteroidetes	251	113	65	43	32	29	9
	Tenericutes	99	90	33	25	13	14	0
	Cyanobacteria	289	112	73	52	39	30	8
	Spirochaetes	171	104	56	41	24	25	5
	Thermotogae	231	110	60	48	36	22	4
	Rest of Bacteria *	255	113	67	48	37	27	6
	PVC	206	110	58	43	28	27	6
Eukarya	Fungi	298	127	105	87	51	52	10
	Metazoa	307	135	136	126	65	75	42
	Plants	332	145	117	87	58	54	14
	Protista	220	117	94	67	39	46	9

* Includes proteomes from Chlorobi, Chloroflexi, Aquificae, Deinococcus thermus, Fusobacteria, Acidobacteria, Deferribacters, Dictyoglomi, Elusimicrobia, Synergistetes, Fibrobacters, Gemmatimonadetes, Nitrospirae, and Thermobaculum.

Table S2. Average percentage of FSF domains in each phyla/kingdom corresponding to the seven general functional categories. Numbers were rounded up when the decimal value exceeded 0.5 and rounded down otherwise. Nanoarchaeota (highlighted in bold) is an outlier considering it has the smallest percentage for metabolic domains compared to the rest and this decrease is offset by an increase in the informational FSFs.

Superkingdom	Phyla/Kingdom	Metabolism	Information	ICP	Regulation	Other	General	ECP
Archaea	Crenarchaeota	48	21	10	9	7	5	1
	Euryarchaeota	47	20	11	9	7	5	1
	Korarchaeota	46	22	10	9	7	5	1
	Nanoarchaeota	29	38	12	8	8	6	1
	Thaumarchaeota	46	21	11	10	5	6	1
Bacteria	Proteobacteria	45	20	13	8	7	5	1
	Firmicutes	44	21	12	10	6	5	1
	Actinobacteria	48	20	12	9	6	5	1
	Bacteroidetes	46	22	12	8	6	5	2
	Tenericutes	36	33	12	9	5	5	0

Table S2. *Cont.*

Superkingdom	Phyla/Kingdom	Metabolism	Information	ICP	Regulation	Other	General	ECP
Bacteria	Cyanobacteria	48	19	12	9	6	5	1
	Spirochaetes	39	25	13	10	6	6	1
	Thermotogae	45	22	12	9	7	4	1
	Rest of Bacteria *	46	21	12	9	7	5	1
	PVC	42	24	12	9	6	6	1
Eukarya	Fungi	41	17	14	12	7	7	1
	Metazoa	35	15	15	14	7	8	5
	Plants	41	18	14	11	7	7	2
	Protista	36	20	16	11	6	8	2

* Includes proteomes from Chlorobi, Chloroflexi, Aquificae, Deinococcus thermus, Fusobacteria, Acidobacteria, Deferribacters, Dictyoglomi, Elusimicrobia, Synergistetes, Fibrobacters, Gemmatimonadetes, Nitrospirae, and Thermobaculum

Table S3. Comparison of functional categories across superkingdoms using Welch's ANOVA.

Functional category	F-ratio	DF	P-value *
<i>Metabolism</i>	350.21	2	<0.0001
<i>Information</i>	582.28	2	<0.0001
<i>ICP</i>	1271.32	2	<0.0001
<i>Regulation</i>	966.75	2	<0.0001
<i>Other</i>	520.97	2	<0.0001
<i>General</i>	1043.76	2	<0.0001
<i>ECP</i>	263.44	2	<0.0001

* All the P-values are statistically significant at 0.05.

Table S4. Names and description of FSF domains corresponding to subcategory *structural proteins* in the main category *General*.

No.	SCOP Id	FSF Id	Description
1	103589	g.71.1	Mini-collagen I, C-terminal domain
2	49695	b.11.1	Gamma-crystallin-like
3	51269	b.85.1	Anti-freeze protein (AFP) III-like domain
4	56558	d.182.1	Baseplate structural protein gp11
5	58002	h.1.6	Chicken cartilage matrix protein
6	58006	h.1.7	Assembly domain of cartilage oligomeric matrix protein
7	75404	d.213.1	Vesiculovirus (VSV) matrix proteins

Table S5. List of organisms analyzed with their taxonomic classifications.

No.	Genome Name	Phyla/Kingdom	Superkingdom
1	<i>Malassezia globosa</i> CBS 7966	Fungi	Eukaryota
2	<i>Ustilago maydis</i>	Fungi	Eukaryota
3	<i>Puccinia graminis f. sp. tritici</i> CRL 75-36-700-3	Fungi	Eukaryota
4	<i>Melampsora laricis-populina</i>	Fungi	Eukaryota
5	<i>Sporobolomyces roseus</i> IAM 13481	Fungi	Eukaryota

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
6	<i>Serpula lacrymans</i> var. <i>lacrymans</i> S7.9	Fungi	Eukaryota
7	<i>Coprinopsis cinerea</i> okayama7 130 v3	Fungi	Eukaryota
8	<i>Pleurotus ostreatus</i>	Fungi	Eukaryota
9	<i>Laccaria bicolor</i> S238N-H82	Fungi	Eukaryota
10	<i>Agaricus bisporus</i> var. <i>bisporus</i>	Fungi	Eukaryota
11	<i>Schizophyllum commune</i>	Fungi	Eukaryota
12	<i>Heterobasidion annosum</i>	Fungi	Eukaryota
13	<i>Phanerochaete chrysosporium</i> RP-78 2.1	Fungi	Eukaryota
14	<i>Postia placenta</i>	Fungi	Eukaryota
15	<i>Tremella mesenterica</i>	Fungi	Eukaryota
16	<i>Cryptococcus neoformans</i> JEC21	Fungi	Eukaryota
17	<i>Magnaporthe grisea</i> 70-15	Fungi	Eukaryota
18	<i>Podospora anserina</i>	Fungi	Eukaryota
19	<i>Sporotrichum thermophile</i> ATCC 42464	Fungi	Eukaryota
20	<i>Thielavia terrestris</i> NRRL 8126	Fungi	Eukaryota
21	<i>Chaetomium globosum</i> CBS 148.51	Fungi	Eukaryota
22	<i>Neurospora tetrasperma</i>	Fungi	Eukaryota
23	<i>Neurospora discreta</i> FGSC 8579	Fungi	Eukaryota
24	<i>Neurospora crassa</i> OR74A	Fungi	Eukaryota
25	<i>Cryphonectria parasitica</i>	Fungi	Eukaryota
26	<i>Verticillium dahliae</i> VdLs.17	Fungi	Eukaryota
27	<i>Verticillium albo-atrum</i> VaMs.102	Fungi	Eukaryota
28	<i>Fusarium oxysporum</i> f. sp. <i>lycopersici</i> 4286	Fungi	Eukaryota
29	<i>Nectria haematococca</i> mpVI	Fungi	Eukaryota
30	<i>Fusarium verticillioides</i> 7600	Fungi	Eukaryota
31	<i>Fusarium graminearum</i>	Fungi	Eukaryota
32	<i>Trichoderma atroviride</i>	Fungi	Eukaryota
33	<i>Trichoderma reesei</i> 1.2	Fungi	Eukaryota
34	<i>Trichoderma virens</i> Gv29-8	Fungi	Eukaryota
35	<i>Botrytis cinerea</i> B05.10	Fungi	Eukaryota
36	<i>Sclerotinia sclerotiorum</i>	Fungi	Eukaryota
37	<i>Alternaria brassicicola</i>	Fungi	Eukaryota
38	<i>Pyrenophora tritici-repentis</i>	Fungi	Eukaryota
39	<i>Cochliobolus heterostrophus</i>	Fungi	Eukaryota
40	<i>Stagonospora nodorum</i>	Fungi	Eukaryota
41	<i>Mycosphaerella fijiensis</i> CIRAD86	Fungi	Eukaryota
42	<i>Mycosphaerella graminicola</i> IPO323	Fungi	Eukaryota
43	<i>Ajellomyces dermatitidis</i> SLH14081	Fungi	Eukaryota
44	<i>Histoplasma capsulatum</i> class NAml strain WU24	Fungi	Eukaryota
45	<i>Microsporium canis</i> CBS 113480	Fungi	Eukaryota
46	<i>Microsporium gypseum</i>	Fungi	Eukaryota
47	<i>Arthroderma benhamiae</i> CBS 112371	Fungi	Eukaryota
48	<i>Trichophyton equinum</i> CBS 127.97	Fungi	Eukaryota
49	<i>Trichophyton verrucosum</i> HKI 0517	Fungi	Eukaryota

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
50	<i>Trichophyton tonsurans</i> CBS 112818	Fungi	Eukaryota
51	<i>Trichophyton rubrum</i> CBS 118892	Fungi	Eukaryota
52	<i>Paracoccidioides brasiliensis</i> Pb18	Fungi	Eukaryota
53	<i>Coccidioides posadasii</i> RMSCC 3488	Fungi	Eukaryota
54	<i>Coccidioides immitis</i> RS	Fungi	Eukaryota
55	<i>Uncinocarpus reesii</i> 1704	Fungi	Eukaryota
56	<i>Aspergillus fumigatus</i> Af293	Fungi	Eukaryota
57	<i>Neosartorya fischeri</i> NRRL 181	Fungi	Eukaryota
58	<i>Penicillium chrysogenum</i> Wisconsin 54-1255	Fungi	Eukaryota
59	<i>Penicillium marneffei</i> ATCC 18224	Fungi	Eukaryota
60	<i>Aspergillus carbonarius</i> ITEM 5010	Fungi	Eukaryota
61	<i>Aspergillus terreus</i> NIH2624	Fungi	Eukaryota
62	<i>Aspergillus oryzae</i> RIB40	Fungi	Eukaryota
63	<i>Aspergillus niger</i> ATCC 1015	Fungi	Eukaryota
64	<i>Aspergillus flavus</i> NRRL3357	Fungi	Eukaryota
65	<i>Aspergillus clavatus</i> NRRL 1	Fungi	Eukaryota
66	<i>Aspergillus nidulans</i> FGSC A4	Fungi	Eukaryota
67	<i>Tuber melanosporum</i> Vittad	Fungi	Eukaryota
68	<i>Pichia stipitis</i> CBS 6054	Fungi	Eukaryota
69	<i>Candida guilliermondii</i> ATCC 6260	Fungi	Eukaryota
70	<i>Lodderomyces elongisporus</i> NRRL YB-4239	Fungi	Eukaryota
71	<i>Debaromyces hansenii</i>	Fungi	Eukaryota
72	<i>Candida dubliniensis</i> CD36	Fungi	Eukaryota
73	<i>Candida tropicalis</i> MYA-3404	Fungi	Eukaryota
74	<i>Candida parapsilosis</i>	Fungi	Eukaryota
75	<i>Candida albicans</i> SC5314	Fungi	Eukaryota
76	<i>Yarrowia lipolytica</i> CLIB122	Fungi	Eukaryota
77	<i>Candida lusitaniae</i> ATCC 42720	Fungi	Eukaryota
78	<i>Vanderwaltozyma polyspora</i> DSM 70294	Fungi	Eukaryota
79	<i>Candida glabrata</i> CBS138	Fungi	Eukaryota
80	<i>Kluyveromyces thermotolerans</i> CBS 6340	Fungi	Eukaryota
81	<i>Lachancea kluyveri</i>	Fungi	Eukaryota
82	<i>Kluyveromyces waltii</i>	Fungi	Eukaryota
83	<i>Ashbya gossypii</i> ATCC 10895	Fungi	Eukaryota
84	<i>Zygosaccharomyces rouxii</i>	Fungi	Eukaryota
85	<i>Saccharomyces mikatae</i> MIT	Fungi	Eukaryota
86	<i>Saccharomyces paradoxus</i> MIT	Fungi	Eukaryota
87	<i>Saccharomyces cerevisiae</i> SGD	Fungi	Eukaryota
88	<i>Saccharomyces bayanus</i> MIT	Fungi	Eukaryota
89	<i>Pichia pastoris</i> GS115	Fungi	Eukaryota
90	<i>Kluyveromyces lactis</i>	Fungi	Eukaryota
91	<i>Schizosaccharomyces octosporus</i> yFS286	Fungi	Eukaryota
92	<i>Schizosaccharomyces japonicus</i> yFS275	Fungi	Eukaryota
93	<i>Schizosaccharomyces pombe</i>	Fungi	Eukaryota

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
94	<i>Allomyces macrogynus</i> ATCC 38327	Fungi	Eukaryota
95	<i>Rhizopus oryzae</i> RA 99-880	Fungi	Eukaryota
96	<i>Phycomyces blakesleeanus</i>	Fungi	Eukaryota
97	<i>Mucor circinelloides</i>	Fungi	Eukaryota
98	<i>Spizellomyces punctatus</i> DAOM BR117	Fungi	Eukaryota
99	<i>Batrachochytrium dendrobatidis</i> JEL423	Fungi	Eukaryota
100	<i>Encephalitozoon cuniculi</i>	Fungi	Eukaryota
101	<i>Encephalitozoon intestinalis</i>	Fungi	Eukaryota
102	<i>Homo sapiens</i> 59_37d (all transcripts)	Metazoa	Eukaryota
103	<i>Pan troglodytes</i> 59_21n (all transcripts)	Metazoa	Eukaryota
104	<i>Gorilla gorilla</i> 59_3b (all transcripts)	Metazoa	Eukaryota
105	<i>Pongo pygmaeus</i> 59_1e (all transcripts)	Metazoa	Eukaryota
106	<i>Macaca mulatta</i> 59_10n (all transcripts)	Metazoa	Eukaryota
107	<i>Callithrix jacchus</i> 59_321a (all transcripts)	Metazoa	Eukaryota
108	<i>Otolemur garnettii</i> 59_1g (all transcripts)	Metazoa	Eukaryota
109	<i>Microcebus murinus</i> 59_1d (all transcripts)	Metazoa	Eukaryota
110	<i>Tarsius syrichta</i> 59_1e (all transcripts)	Metazoa	Eukaryota
111	<i>Rattus norvegicus</i> 59_34a (all transcripts)	Metazoa	Eukaryota
112	<i>Mus musculus</i> 59_37l (all transcripts)	Metazoa	Eukaryota
113	<i>Spermophilus tridecemlineatus</i> 59_1i (all transcripts)	Metazoa	Eukaryota
114	<i>Dipodomys ordii</i> 59_1e (all transcripts)	Metazoa	Eukaryota
115	<i>Cavia porcellus</i> 59_3c (all transcripts)	Metazoa	Eukaryota
116	<i>Oryctolagus cuniculus</i> 59_2b (all transcripts)	Metazoa	Eukaryota
117	<i>Ochotona princeps</i> 59_1e (all transcripts)	Metazoa	Eukaryota
118	<i>Tupaia belangeri</i> 59_1h (all transcripts)	Metazoa	Eukaryota
119	<i>Sus scrofa</i> 59_9c (all transcripts)	Metazoa	Eukaryota
120	<i>Bos taurus</i> 59_4h (all transcripts)	Metazoa	Eukaryota
121	<i>Vicugna pacos</i> 59_1e (all transcripts)	Metazoa	Eukaryota
122	<i>Tursiops truncatus</i> 59_1e (all transcripts)	Metazoa	Eukaryota
123	<i>Canis familiaris</i> 59_2o (all transcripts)	Metazoa	Eukaryota
124	<i>Felis catus</i> 59_1h (all transcripts)	Metazoa	Eukaryota
125	<i>Equus caballus</i> 59_2f (all transcripts)	Metazoa	Eukaryota
126	<i>Myotis lucifugus</i> 59_1i (all transcripts)	Metazoa	Eukaryota
127	<i>Pteropus vampyrus</i> 59_1e (all transcripts)	Metazoa	Eukaryota
128	<i>Sorex araneus</i> 59_1g (all transcripts)	Metazoa	Eukaryota
129	<i>Erinaceus europaeus</i> 59_1g (all transcripts)	Metazoa	Eukaryota
130	<i>Procapra capensis</i> 59_1e (all transcripts)	Metazoa	Eukaryota
131	<i>Loxodonta africana</i> 59_3b (all transcripts)	Metazoa	Eukaryota
132	<i>Echinops telfairi</i> 59_1i (all transcripts)	Metazoa	Eukaryota
133	<i>Dasyurus novemcinctus</i> 59_2c (all transcripts)	Metazoa	Eukaryota
134	<i>Macropus eugenii</i> 59_1b (all transcripts)	Metazoa	Eukaryota
135	<i>Monodelphis domestica</i> 59_5k (all transcripts)	Metazoa	Eukaryota
136	<i>Ornithorhynchus anatinus</i> 59_1m (all transcripts)	Metazoa	Eukaryota
137	<i>Anolis carolinensis</i> 59_1c (all transcripts)	Metazoa	Eukaryota

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
138	<i>Taeniopygia guttata</i> 59_1e (all transcripts)	Metazoa	Eukaryota
139	<i>Meleagris gallopavo</i> 57_2 (all transcripts)	Metazoa	Eukaryota
140	<i>Gallus gallus</i> 59_2o (all transcripts)	Metazoa	Eukaryota
141	<i>Xenopus laevis</i>	Metazoa	Eukaryota
142	<i>Xenopus tropicalis</i> 59_41p (all transcripts)	Metazoa	Eukaryota
143	<i>Danio rerio</i> 59_8e (all transcripts)	Metazoa	Eukaryota
144	<i>Gasterosteus aculeatus</i> 59_1l (all transcripts)	Metazoa	Eukaryota
145	<i>Oryzias latipes</i> 59_1k (all transcripts)	Metazoa	Eukaryota
146	<i>Tetraodon nigroviridis</i> 59_8d (all transcripts)	Metazoa	Eukaryota
147	<i>Takifugu rubripes</i> 59_4m (all transcripts)	Metazoa	Eukaryota
148	<i>Branchiostoma floridae</i> 1.0	Metazoa	Eukaryota
149	<i>Ciona savignyi</i> 59_2j (all transcripts)	Metazoa	Eukaryota
150	<i>Ciona intestinalis</i> 59_2o (all transcripts)	Metazoa	Eukaryota
151	<i>Strongylocentrotus purpuratus</i>	Metazoa	Eukaryota
152	<i>Helobdella robusta</i>	Metazoa	Eukaryota
153	<i>Capitella</i> sp. I	Metazoa	Eukaryota
154	<i>Bombyx mori</i>	Metazoa	Eukaryota
155	<i>Nasonia vitripennis</i>	Metazoa	Eukaryota
156	<i>Apis mellifera</i> 38.2d (all transcripts)	Metazoa	Eukaryota
157	<i>Drosophila grimshawi</i> 1.3	Metazoa	Eukaryota
158	<i>Drosophila willistoni</i> 1.3	Metazoa	Eukaryota
159	<i>Drosophila pseudoobscura</i> 2.13	Metazoa	Eukaryota
160	<i>Drosophila persimilis</i> 1.3	Metazoa	Eukaryota
161	<i>Drosophila yakuba</i> 1.3	Metazoa	Eukaryota
162	<i>Drosophila simulans</i> 1.3	Metazoa	Eukaryota
163	<i>Drosophila sechellia</i> 1.3	Metazoa	Eukaryota
164	<i>Drosophila melanogaster</i> 59_525a (all transcripts)	Metazoa	Eukaryota
165	<i>Drosophila erecta</i> 1.3	Metazoa	Eukaryota
166	<i>Drosophila ananassae</i> 1.3	Metazoa	Eukaryota
167	<i>Drosophila virilis</i> 1.2	Metazoa	Eukaryota
168	<i>Drosophila mojavensis</i> 1.3	Metazoa	Eukaryota
169	<i>Aedes aegypti</i> 55 (all transcripts)	Metazoa	Eukaryota
170	<i>Culex pipiens quinquefasciatus</i>	Metazoa	Eukaryota
171	<i>Anopheles gambiae</i> 49_3j (all transcripts)	Metazoa	Eukaryota
172	<i>Tribolium castaneum</i> 3.0	Metazoa	Eukaryota
173	<i>Pediculus humanus corporis</i>	Metazoa	Eukaryota
174	<i>Acyrtosiphon pisum</i>	Metazoa	Eukaryota
175	<i>Daphnia pulex</i>	Metazoa	Eukaryota
176	<i>Ixodes scapularis</i>	Metazoa	Eukaryota
177	<i>Lottia gigantea</i>	Metazoa	Eukaryota
178	<i>Pristionchus pacificus</i>	Metazoa	Eukaryota
179	<i>Meloidogyne incognita</i>	Metazoa	Eukaryota
180	<i>Brugia malayi</i> WS218	Metazoa	Eukaryota
181	<i>Caenorhabditis japonica</i>	Metazoa	Eukaryota

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
182	<i>Caenorhabditis brenneri</i>	Metazoa	Eukaryota
183	<i>Caenorhabditis remanei</i>	Metazoa	Eukaryota
184	<i>Caenorhabditis elegans</i> 59_210a (all transcripts)	Metazoa	Eukaryota
185	<i>Caenorhabditis briggsae</i> 2	Metazoa	Eukaryota
186	<i>Schistosoma mansoni</i>	Metazoa	Eukaryota
187	<i>Nematostella vectensis</i> 1.0	Metazoa	Eukaryota
188	<i>Hydra magnipapillata</i>	Metazoa	Eukaryota
189	<i>Trichoplax adhaerens</i>	Metazoa	Eukaryota
190	<i>Giardia lamblia</i> 2.3	Protista	Eukaryota
191	<i>Trypanosoma cruzi</i> strain CL Brener	Protista	Eukaryota
192	<i>Trypanosoma brucei</i>	Protista	Eukaryota
193	<i>Leishmania mexicana</i> 2.4	Protista	Eukaryota
194	<i>Leishmania major</i> strain Friedlin	Protista	Eukaryota
195	<i>Leishmania infantum</i> JPCM5 2.4	Protista	Eukaryota
196	<i>Leishmania braziliensis</i> MHOM/BR/75/M2904 2.4	Protista	Eukaryota
197	<i>Aureococcus anophagefferens</i>	Protista	Eukaryota
198	<i>Phytophthora ramorum</i> 1.1	Protista	Eukaryota
199	<i>Phytophthora sojae</i> 1.1	Protista	Eukaryota
200	<i>Phytophthora infestans</i> T30-4	Protista	Eukaryota
201	<i>Phytophthora capsici</i>	Protista	Eukaryota
202	<i>Paramecium tetraurelia</i>	Protista	Eukaryota
203	<i>Tetrahymena thermophila</i> SB210 1	Protista	Eukaryota
204	<i>Babesia bovis</i> T2Bo	Protista	Eukaryota
205	<i>Theileria parva</i>	Protista	Eukaryota
206	<i>Theileria annulata</i>	Protista	Eukaryota
207	<i>Plasmodium falciparum</i> 3D7	Protista	Eukaryota
208	<i>Plasmodium vivax</i> SaI-1 7.0	Protista	Eukaryota
209	<i>Plasmodium knowlesi</i> strain H	Protista	Eukaryota
210	<i>Plasmodium yoelii</i> ssp. yoelii 1	Protista	Eukaryota
211	<i>Plasmodium chabaudi</i>	Protista	Eukaryota
212	<i>Plasmodium berghei</i> ANKA	Protista	Eukaryota
213	<i>Cryptosporidium hominis</i>	Protista	Eukaryota
214	<i>Cryptosporidium muris</i>	Protista	Eukaryota
215	<i>Cryptosporidium parvum</i> Iowa II	Protista	Eukaryota
216	<i>Neospora caninum</i> Nc-Liverpool 6.2	Protista	Eukaryota
217	<i>Neospora caninum</i>	Protista	Eukaryota
218	<i>Toxoplasma gondii</i> ME49	Protista	Eukaryota
219	<i>Naegleria gruberi</i>	Protista	Eukaryota
220	<i>Guillardia theta</i>	Protista	Eukaryota
221	<i>Arabidopsis lyrata</i>	Plantae	Eukaryota
222	<i>Arabidopsis thaliana</i> 10 (all transcripts)	Plantae	Eukaryota
223	<i>Carica papaya</i>	Plantae	Eukaryota
224	<i>Medicago truncatula</i>	Plantae	Eukaryota
225	<i>Glycine max</i>	Plantae	Eukaryota

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
226	<i>Cucumis sativus</i>	Plantae	Eukaryota
227	<i>Populus trichocarpa</i> 6.0	Plantae	Eukaryota
228	<i>Vitis vinifera</i>	Plantae	Eukaryota
229	<i>Brachypodium distachyon</i>	Plantae	Eukaryota
230	<i>Oryza sativa</i> ssp. <i>japonica</i> 5.0	Plantae	Eukaryota
231	<i>Zea mays</i> subsp. <i>mays</i>	Plantae	Eukaryota
232	<i>Sorghum bicolor</i>	Plantae	Eukaryota
233	<i>Selaginella moellendorffii</i>	Plantae	Eukaryota
234	<i>Physcomitrella patens</i> subsp. <i>patens</i>	Plantae	Eukaryota
235	<i>Ostreococcus</i> sp. <i>RCC809</i>	Plantae	Eukaryota
236	<i>Ostreococcus lucimarinus</i> <i>CCE9901</i>	Plantae	Eukaryota
237	<i>Ostreococcus tauri</i>	Plantae	Eukaryota
238	<i>Micromonas</i> sp. <i>RCC299</i>	Plantae	Eukaryota
239	<i>Micromonas pusilla</i> <i>CCMP1545</i>	Plantae	Eukaryota
240	<i>Coccomyxa</i> sp. <i>C-169</i>	Plantae	Eukaryota
241	<i>Chlorella</i> sp. <i>NC64A</i>	Plantae	Eukaryota
242	<i>Chlorella vulgaris</i>	Plantae	Eukaryota
243	<i>Volvox carteri</i> f. <i>nagariensis</i>	Plantae	Eukaryota
244	<i>Chlamydomonas reinhardtii</i> 4.0	Plantae	Eukaryota
245	<i>Candidatus Koribacter versatilis</i> <i>Ellin345</i>	Acidobacteria	Bacteria
246	<i>Candidatus Solibacter usitatus</i> <i>Ellin6076</i>	Acidobacteria	Bacteria
247	<i>Acidobacterium capsulatum</i> <i>ATCC 51196</i>	Acidobacteria	Bacteria
248	<i>Gardnerella vaginalis</i> 409-05	Actinobacteria	Bacteria
249	<i>Bifidobacterium longum</i> <i>NCC2705</i>	Actinobacteria	Bacteria
250	<i>Bifidobacterium animalis</i> ssp. <i>lactis</i> <i>AD011</i>	Actinobacteria	Bacteria
251	<i>Bifidobacterium dentium</i> <i>Bd1</i>	Actinobacteria	Bacteria
252	<i>Bifidobacterium adolescentis</i> <i>ATCC 15703</i>	Actinobacteria	Bacteria
253	<i>Kineococcus radiotolerans</i> <i>SRS30216</i>	Actinobacteria	Bacteria
254	<i>Catenulispora acidiphila</i> <i>DSM 44928</i>	Actinobacteria	Bacteria
255	<i>Stackebrandtia nassauensis</i> <i>DSM 44728</i>	Actinobacteria	Bacteria
256	<i>Acidothermus cellulolyticus</i> <i>11B</i>	Actinobacteria	Bacteria
257	<i>Nakamurella multipartita</i> <i>DSM 44233</i>	Actinobacteria	Bacteria
258	<i>Geodermatophilus obscurus</i> <i>DSM 43160</i>	Actinobacteria	Bacteria
259	<i>Frankia</i> sp. <i>Cc13</i>	Actinobacteria	Bacteria
260	<i>Frankia alni</i> <i>ACN14a</i>	Actinobacteria	Bacteria
261	<i>Thermobifida fusca</i> <i>YX</i>	Actinobacteria	Bacteria
262	<i>Thermomonospora curvata</i> <i>DSM 43183</i>	Actinobacteria	Bacteria
263	<i>Streptosporangium roseum</i> <i>DSM 43021</i>	Actinobacteria	Bacteria
264	<i>Streptomyces griseus</i> ssp. <i>griseus</i> <i>NBRC 13350</i>	Actinobacteria	Bacteria
265	<i>Streptomyces avermitilis</i> <i>MA-4680</i>	Actinobacteria	Bacteria
266	<i>Streptomyces scabiei</i> 87.22	Actinobacteria	Bacteria
267	<i>Streptomyces coelicolor</i>	Actinobacteria	Bacteria
268	<i>Actinosynnema mirum</i> <i>DSM 43827</i>	Actinobacteria	Bacteria
269	<i>Saccharomonospora viridis</i> <i>DSM 43017</i>	Actinobacteria	Bacteria

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
270	<i>Saccharopolyspora erythraea</i> NRRL 2338	Actinobacteria	Bacteria
271	<i>Kribbella flavida</i> DSM 17836	Actinobacteria	Bacteria
272	<i>Nocardioides</i> sp. JS614	Actinobacteria	Bacteria
273	<i>Propionibacterium acnes</i> KPA171202	Actinobacteria	Bacteria
274	<i>Salinispora arenicola</i> CNS-205	Actinobacteria	Bacteria
275	<i>Salinispora tropica</i> CNB-440	Actinobacteria	Bacteria
276	<i>Gordonia bronchialis</i> DSM 43247	Actinobacteria	Bacteria
277	<i>Rhodococcus jostii</i> RHA1	Actinobacteria	Bacteria
278	<i>Rhodococcus opacus</i> B4	Actinobacteria	Bacteria
279	<i>Rhodococcus erythropolis</i> PR4	Actinobacteria	Bacteria
280	<i>Nocardia farcinica</i> IFM 10152	Actinobacteria	Bacteria
281	<i>Mycobacterium abscessus</i> ATCC 19977	Actinobacteria	Bacteria
282	<i>Mycobacterium</i> sp. MCS	Actinobacteria	Bacteria
283	<i>Mycobacterium avium</i> ssp. <i>paratuberculosis</i> K-10	Actinobacteria	Bacteria
284	<i>Mycobacterium vanbaalenii</i> PYR-1	Actinobacteria	Bacteria
285	<i>Mycobacterium tuberculosis</i> H37Rv	Actinobacteria	Bacteria
286	<i>Mycobacterium bovis</i> AF2122/97	Actinobacteria	Bacteria
287	<i>Mycobacterium ulcerans</i> Agy99	Actinobacteria	Bacteria
288	<i>Mycobacterium gilvum</i> PYR-GCK	Actinobacteria	Bacteria
289	<i>Mycobacterium marinum</i> M	Actinobacteria	Bacteria
290	<i>Mycobacterium smegmatis</i> MC2 155	Actinobacteria	Bacteria
291	<i>Mycobacterium leprae</i> TN	Actinobacteria	Bacteria
292	<i>Corynebacterium aurimucosum</i> ATCC 700975	Actinobacteria	Bacteria
293	<i>Corynebacterium kroppenstedtii</i> DSM 44385	Actinobacteria	Bacteria
294	<i>Corynebacterium efficiens</i> YS-314	Actinobacteria	Bacteria
295	<i>Corynebacterium urealyticum</i> DSM 7109	Actinobacteria	Bacteria
296	<i>Corynebacterium jeikeium</i> K411	Actinobacteria	Bacteria
297	<i>Corynebacterium glutamicum</i> ATCC 13032 Kitasato	Actinobacteria	Bacteria
298	<i>Corynebacterium diphtheriae</i> NCTC 13129	Actinobacteria	Bacteria
299	<i>Tropheryma whipplei</i> Twist	Actinobacteria	Bacteria
300	<i>Sanguibacter keddieii</i> DSM 10542	Actinobacteria	Bacteria
301	<i>Kytococcus sedentarius</i> DSM 20547	Actinobacteria	Bacteria
302	<i>Beutenbergia cavernae</i> DSM 12333	Actinobacteria	Bacteria
303	<i>Leifsonia xyli</i> ssp. <i>xyli</i> CTCB07	Actinobacteria	Bacteria
304	<i>Clavibacter michiganensis</i> ssp. <i>michiganensis</i> NCPPB 382	Actinobacteria	Bacteria
305	<i>Jonesia denitrificans</i> DSM 20603	Actinobacteria	Bacteria
306	<i>Brachybacterium faecium</i> DSM 4810	Actinobacteria	Bacteria
307	<i>Xylanimonas cellulosilytica</i> DSM 15894	Actinobacteria	Bacteria
308	<i>Kocuria rhizophila</i> DC2201	Actinobacteria	Bacteria
309	<i>Rothia mucilaginosa</i> DY-18	Actinobacteria	Bacteria
310	<i>Arthrobacter</i> sp. FB24	Actinobacteria	Bacteria
311	<i>Arthrobacter chlorophenolicus</i> A6	Actinobacteria	Bacteria
312	<i>Arthrobacter aurescens</i> TC1	Actinobacteria	Bacteria
313	<i>Renibacterium salmoninarum</i> ATCC 33209	Actinobacteria	Bacteria

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
314	<i>Micrococcus luteus</i> NCTC 2665	Actinobacteria	Bacteria
315	<i>Cryptobacterium curtum</i> DSM 15641	Actinobacteria	Bacteria
316	<i>Eggerthella lenta</i> DSM 2243	Actinobacteria	Bacteria
317	<i>Slackia heliotrinireducens</i> DSM 20476	Actinobacteria	Bacteria
318	<i>Atopobium parvulum</i> DSM 20469	Actinobacteria	Bacteria
319	<i>Conexibacter woesei</i> DSM 14684	Actinobacteria	Bacteria
320	<i>Rubrobacter xylanophilus</i> DSM 9941	Actinobacteria	Bacteria
321	<i>Acidimicrobium ferrooxidans</i> DSM 10331	Actinobacteria	Bacteria
322	<i>Sulfurihydrogenibium</i> sp. YO3AOP1	Aquificae	Bacteria
323	<i>Sulfurihydrogenibium azorense</i> Az-Fu1	Aquificae	Bacteria
324	<i>Persephonella marina</i> EX-H1	Aquificae	Bacteria
325	<i>Hydrogenobaculum</i> sp. Y04AAS1	Aquificae	Bacteria
326	<i>Thermocrinis albus</i> DSM 14484	Aquificae	Bacteria
327	<i>Aquifex aeolicus</i> VF5	Aquificae	Bacteria
328	<i>Hydrogenobacter thermophilus</i> TK-6	Aquificae	Bacteria
329	<i>Dyadobacter fermentans</i> DSM 18053	Bacteroidetes	Bacteria
330	<i>Cytophaga hutchinsonii</i> ATCC 33406	Bacteroidetes	Bacteria
331	<i>Spirosoma linguale</i> DSM 74	Bacteroidetes	Bacteria
332	<i>Candidatus Azobacteroides pseudotrichonymphae</i> genomovar.	Bacteroidetes	Bacteria
333	<i>Prevotella ruminicola</i> 23	Bacteroidetes	Bacteria
334	<i>Parabacteroides distasonis</i> ATCC 8503	Bacteroidetes	Bacteria
335	<i>Porphyromonas gingivalis</i> W83	Bacteroidetes	Bacteria
336	<i>Bacteroides vulgatus</i> ATCC 8482	Bacteroidetes	Bacteria
337	<i>Bacteroides thetaiotaomicron</i> VPI-5482	Bacteroidetes	Bacteria
338	<i>Bacteroides fragilis</i> NCTC 9343	Bacteroidetes	Bacteria
339	<i>Candidatus Amoebophilus asiaticus</i> 5a2	Bacteroidetes	Bacteria
340	<i>Salinibacter ruber</i> DSM 13855	Bacteroidetes	Bacteria
341	<i>Rhodothermus marinus</i> DSM 4252	Bacteroidetes	Bacteria
342	<i>Chitinophaga pinensis</i> DSM 2588	Bacteroidetes	Bacteria
343	<i>Pedobacter heparinus</i> DSM 2366	Bacteroidetes	Bacteria
344	<i>Candidatus Sulcia muelleri</i> GWSS	Bacteroidetes	Bacteria
345	<i>Zunongwangia profunda</i> SM-A87	Bacteroidetes	Bacteria
346	<i>Gramella forsetii</i> KT0803	Bacteroidetes	Bacteria
347	<i>Robiginitalea biformata</i> HTCC2501	Bacteroidetes	Bacteria
348	<i>Flavobacteriaceae bacterium</i> 3519-10	Bacteroidetes	Bacteria
349	<i>Capnocytophaga ochracea</i> DSM 7271	Bacteroidetes	Bacteria
350	<i>Flavobacterium psychrophilum</i> JIP02/86	Bacteroidetes	Bacteria
351	<i>Flavobacterium johnsoniae</i> UW101	Bacteroidetes	Bacteria
352	<i>Blattabacterium</i> sp. Bge	Bacteroidetes	Bacteria
353	<i>Candidatus Protochlamydia amoebophila</i> UWE25	Chlamydiae	Bacteria
354	<i>Chlamydophila pneumoniae</i> TW-183	Chlamydiae	Bacteria
355	<i>Chlamydophila caviae</i> GPIC	Chlamydiae	Bacteria
356	<i>Chlamydophila felis</i> Fe/C-56	Chlamydiae	Bacteria
357	<i>Chlamydophila abortus</i> S26/3	Chlamydiae	Bacteria

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
358	<i>Chlamydia muridarum</i> Nigg	Chlamydiae	Bacteria
359	<i>Chlamydia trachomatis</i> D/UW-3/CX	Chlamydiae	Bacteria
360	<i>Pelodictyon phaeoclathratiforme</i> BU-1	Chlorobi	Bacteria
361	<i>Chlorobium luteolum</i> DSM 273	Chlorobi	Bacteria
362	<i>Chlorobium chlorochromatii</i> CaD3	Chlorobi	Bacteria
363	<i>Chlorobium phaeobacteroides</i> DSM 266	Chlorobi	Bacteria
364	<i>Chlorobium phaeovibrioides</i> DSM 265	Chlorobi	Bacteria
365	<i>Chlorobium limicola</i> DSM 245	Chlorobi	Bacteria
366	<i>Chlorobaculum parvum</i> NCIB 8327	Chlorobi	Bacteria
367	<i>Chlorobium tepidum</i> TLS	Chlorobi	Bacteria
368	<i>Chloroherpeton thalassium</i> ATCC 35110	Chlorobi	Bacteria
369	<i>Prosthecochloris aestuarii</i> DSM 271	Chlorobi	Bacteria
370	<i>Dehalococcoides</i> sp. CBDB1	Chloroflexi	Bacteria
371	<i>Dehalococcoides ethenogenes</i> 195	Chloroflexi	Bacteria
372	<i>Thermomicrobium roseum</i> DSM 5159	Chloroflexi	Bacteria
373	<i>Sphaerobacter thermophilus</i> DSM 20745	Chloroflexi	Bacteria
374	<i>Herpetosiphon aurantiacus</i> ATCC 23779	Chloroflexi	Bacteria
375	<i>Roseiflexus</i> sp. RS-1	Chloroflexi	Bacteria
376	<i>Roseiflexus castenholzii</i> DSM 13941	Chloroflexi	Bacteria
377	<i>Chloroflexus</i> sp. Y-400-fl	Chloroflexi	Bacteria
378	<i>Chloroflexus aggregans</i> DSM 9485	Chloroflexi	Bacteria
379	<i>Chloroflexus aurantiacus</i> J-10-fl	Chloroflexi	Bacteria
380	<i>Gloeobacter violaceus</i> PCC 7421	Cyanobacteria	Bacteria
381	<i>Acaryochloris marina</i> MBIC11017	Cyanobacteria	Bacteria
382	<i>Prochlorococcus marinus</i> MIT 9313	Cyanobacteria	Bacteria
383	<i>Nostoc punctiforme</i> PCC 73102	Cyanobacteria	Bacteria
384	<i>Nostoc</i> sp. PCC 7120	Cyanobacteria	Bacteria
385	<i>Anabaena variabilis</i> ATCC 29413	Cyanobacteria	Bacteria
386	<i>Trichodesmium erythraeum</i> IMS101	Cyanobacteria	Bacteria
387	<i>Thermosynechococcus elongatus</i> BP-1	Cyanobacteria	Bacteria
388	<i>cyanobacterium UCYN-A</i>	Cyanobacteria	Bacteria
389	<i>Cyanothece</i> sp. ATCC 51142	Cyanobacteria	Bacteria
390	<i>Synechocystis</i> sp. PCC 6803	Cyanobacteria	Bacteria
391	<i>Synechococcus elongatus</i> PCC 6301	Cyanobacteria	Bacteria
392	<i>Microcystis aeruginosa</i> NIES-843	Cyanobacteria	Bacteria
393	<i>Denitrovibrio acetiphilus</i> DSM 12809	Deferribacteres	Bacteria
394	<i>Deferribacter desulfuricans</i> SSM1	Deferribacteres	Bacteria
395	<i>Deinococcus deserti</i> VCD115	Deinococcus-Thermus	Bacteria
396	<i>Deinococcus geothermalis</i> DSM 11300	Deinococcus-Thermus	Bacteria
397	<i>Deinococcus radiodurans</i> R1	Deinococcus-Thermus	Bacteria
398	<i>Meiothermus ruber</i> DSM 1279	Deinococcus-Thermus	Bacteria
399	<i>Thermus thermophilus</i> HB27	Deinococcus-Thermus	Bacteria
400	<i>Dictyoglomus turgidum</i> DSM 6724	Dictyoglomi	Bacteria
401	<i>Dictyoglomus thermophilum</i> H-6-12	Dictyoglomi	Bacteria

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
402	<i>Elusimicrobium minutum</i> Pei191	Elusimicrobia	Bacteria
403	uncultured Termite group 1 bacterium phylotype Rs-D17	Elusimicrobia	Bacteria
404	<i>Fibrobacter succinogenes</i> ssp. <i>succinogenes</i> S85	Fibrobacteres	Bacteria
405	<i>Acidaminococcus fermentans</i> DSM 20731	Firmicutes	Bacteria
406	<i>Veillonella parvula</i> DSM 2008	Firmicutes	Bacteria
407	<i>Natranaerobius thermophilus</i> JW/NM-WN-LF	Firmicutes	Bacteria
408	<i>Symbiobacterium thermophilum</i> IAM 14863	Firmicutes	Bacteria
409	<i>Anaerococcus prevotii</i> DSM 20548	Firmicutes	Bacteria
410	<i>Finegoldia magna</i> ATCC 29328	Firmicutes	Bacteria
411	<i>Clostridiales</i> genomosp. BVAB3 UPII9-5	Firmicutes	Bacteria
412	<i>Candidatus Desulfurudis audaxviator</i> MP104C	Firmicutes	Bacteria
413	<i>Pelotomaculum thermopropionicum</i> SI	Firmicutes	Bacteria
414	<i>Desulfitobacterium hafniense</i> Y51	Firmicutes	Bacteria
415	<i>Desulfotomaculum reducens</i> MI-1	Firmicutes	Bacteria
416	<i>Desulfotomaculum acetoxidans</i> DSM 771	Firmicutes	Bacteria
417	<i>Eubacterium rectale</i> ATCC 33656	Firmicutes	Bacteria
418	<i>Eubacterium eligens</i> ATCC 27750	Firmicutes	Bacteria
419	<i>Syntrophomonas wolfei</i> ssp. <i>wolfei</i> Goettingen	Firmicutes	Bacteria
420	<i>Heliobacterium modesticaldum</i> Ice1	Firmicutes	Bacteria
421	<i>Alkaliphilus oremlandii</i> OhILAs	Firmicutes	Bacteria
422	<i>Alkaliphilus metalliredigens</i> QYMF	Firmicutes	Bacteria
423	<i>Clostridium phytofermentans</i> ISDg	Firmicutes	Bacteria
424	<i>Clostridium novyi</i> NT	Firmicutes	Bacteria
425	<i>Clostridium kluyveri</i> DSM 555	Firmicutes	Bacteria
426	<i>Clostridium cellulolyticum</i> H10	Firmicutes	Bacteria
427	<i>Clostridium beijerinckii</i> NCIMB 8052	Firmicutes	Bacteria
428	<i>Clostridium thermocellum</i> ATCC 27405	Firmicutes	Bacteria
429	<i>Clostridium tetani</i> E88	Firmicutes	Bacteria
430	<i>Clostridium perfringens</i> 13	Firmicutes	Bacteria
431	<i>Clostridium difficile</i> 630	Firmicutes	Bacteria
432	<i>Clostridium botulinum</i> A ATCC 3502	Firmicutes	Bacteria
433	<i>Clostridium acetobutylicum</i> ATCC 824	Firmicutes	Bacteria
434	<i>Caldicellulosiruptor saccharolyticus</i> DSM 8903	Firmicutes	Bacteria
435	<i>Anaerocellum thermophilum</i> DSM 6725	Firmicutes	Bacteria
436	<i>Coprothermobacter proteolyticus</i> DSM 5265	Firmicutes	Bacteria
437	<i>Thermoanaerobacter tengcongensis</i> MB4	Firmicutes	Bacteria
438	<i>Carboxydotherrmus hydrogenoformans</i> Z-2901	Firmicutes	Bacteria
439	<i>Moorella thermoacetica</i> ATCC 39073	Firmicutes	Bacteria
440	<i>Ammonifex degensii</i> KC4	Firmicutes	Bacteria
441	<i>Thermoanaerobacter pseudethanolicus</i> ATCC 33223	Firmicutes	Bacteria
442	<i>Thermoanaerobacter</i> sp. X514	Firmicutes	Bacteria
443	<i>Thermoanaerobacter italicus</i> Ab9	Firmicutes	Bacteria
444	<i>Halothermothrix orenii</i> H 168	Firmicutes	Bacteria
445	<i>Enterococcus faecalis</i> V583	Firmicutes	Bacteria

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
446	<i>Oenococcus oeni</i> PSU-1	Firmicutes	Bacteria
447	<i>Leuconostoc citreum</i> KM20	Firmicutes	Bacteria
448	<i>Leuconostoc mesenteroides</i> ssp. <i>mesenteroides</i> ATCC 8293	Firmicutes	Bacteria
449	<i>Lactobacillus casei</i> ATCC 334	Firmicutes	Bacteria
450	<i>Lactobacillus crispatus</i> ST1	Firmicutes	Bacteria
451	<i>Lactobacillus rhamnosus</i> GG	Firmicutes	Bacteria
452	<i>Lactobacillus johnsonii</i> NCC 533	Firmicutes	Bacteria
453	<i>Lactobacillus salivarius</i> UCC118	Firmicutes	Bacteria
454	<i>Lactobacillus fermentum</i> IFO 3956	Firmicutes	Bacteria
455	<i>Lactobacillus sakei</i> ssp. <i>sakei</i> 23K	Firmicutes	Bacteria
456	<i>Lactobacillus reuteri</i> DSM 20016	Firmicutes	Bacteria
457	<i>Lactobacillus gasserii</i> ATCC 33323	Firmicutes	Bacteria
458	<i>Lactobacillus plantarum</i> WCFS1	Firmicutes	Bacteria
459	<i>Lactobacillus helveticus</i> DPC 4571	Firmicutes	Bacteria
460	<i>Lactobacillus delbrueckii</i> ssp. <i>bulgaricus</i> ATCC 11842	Firmicutes	Bacteria
461	<i>Lactobacillus brevis</i> ATCC 367	Firmicutes	Bacteria
462	<i>Lactobacillus acidophilus</i> NCFM	Firmicutes	Bacteria
463	<i>Pediococcus pentosaceus</i> ATCC 25745	Firmicutes	Bacteria
464	<i>Lactococcus lactis</i> ssp. <i>lactis</i> I1403	Firmicutes	Bacteria
465	<i>Streptococcus gallolyticus</i> UCN34	Firmicutes	Bacteria
466	<i>Streptococcus equi</i> ssp. <i>zooepidemicus</i> MGCS10565	Firmicutes	Bacteria
467	<i>Streptococcus dysgalactiae</i> ssp. <i>equisimilis</i> GGS_124	Firmicutes	Bacteria
468	<i>Streptococcus mitis</i> B6	Firmicutes	Bacteria
469	<i>Streptococcus uberis</i> 0140J	Firmicutes	Bacteria
470	<i>Streptococcus pyogenes</i> M1 GAS	Firmicutes	Bacteria
471	<i>Streptococcus pneumoniae</i> TIGR4	Firmicutes	Bacteria
472	<i>Streptococcus agalactiae</i> NEM316	Firmicutes	Bacteria
473	<i>Streptococcus mutans</i> UA159	Firmicutes	Bacteria
474	<i>Streptococcus thermophilus</i> LMG 18311	Firmicutes	Bacteria
475	<i>Streptococcus suis</i> 05ZYH33	Firmicutes	Bacteria
476	<i>Streptococcus sanguinis</i> SK36	Firmicutes	Bacteria
477	<i>Streptococcus gordonii</i> Challis subCH1	Firmicutes	Bacteria
478	<i>Exiguobacterium</i> sp. AT1b	Firmicutes	Bacteria
479	<i>Exiguobacterium sibiricum</i> 255-15	Firmicutes	Bacteria
480	<i>Bacillus tusciae</i> DSM 2912	Firmicutes	Bacteria
481	<i>Alicyclobacillus acidocaldarius</i> ssp. <i>acidocaldarius</i> DSM 446	Firmicutes	Bacteria
482	<i>Brevibacillus brevis</i> NBRC 100599	Firmicutes	Bacteria
483	<i>Paenibacillus</i> sp. JDR-2	Firmicutes	Bacteria
484	<i>Listeria welshimeri</i> ser. 6b SLCC5334	Firmicutes	Bacteria
485	<i>Listeria innocua</i> Clip11262	Firmicutes	Bacteria
486	<i>Listeria seeligeri</i> ser. 1/2b SLCC3954	Firmicutes	Bacteria
487	<i>Listeria monocytogenes</i> EGD-e	Firmicutes	Bacteria
488	<i>Lysinibacillus sphaericus</i> C3-41	Firmicutes	Bacteria
489	<i>Oceanobacillus iheyensis</i> HTE831	Firmicutes	Bacteria

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
490	<i>Anoxybacillus flavithermus</i> WK1	Firmicutes	Bacteria
491	<i>Geobacillus</i> sp. WCH70	Firmicutes	Bacteria
492	<i>Geobacillus thermodenitrificans</i> NG80-2	Firmicutes	Bacteria
493	<i>Geobacillus kaustophilus</i> HTA426	Firmicutes	Bacteria
494	<i>Bacillus subtilis</i> ssp. <i>subtilis</i> 168	Firmicutes	Bacteria
495	<i>Bacillus licheniformis</i> ATCC 14580	Firmicutes	Bacteria
496	<i>Bacillus amyloliquefaciens</i> FZB42	Firmicutes	Bacteria
497	<i>Bacillus halodurans</i> C-125	Firmicutes	Bacteria
498	<i>Bacillus weihenstephanensis</i> KBAB4	Firmicutes	Bacteria
499	<i>Bacillus thuringiensis</i> ser. <i>konkukian</i> 97-27	Firmicutes	Bacteria
500	<i>Bacillus cereus</i> ATCC 14579	Firmicutes	Bacteria
501	<i>Bacillus anthracis</i> Ames Ancestor	Firmicutes	Bacteria
502	<i>Bacillus pseudofirmus</i> OF4	Firmicutes	Bacteria
503	<i>Bacillus clausii</i> KSM-K16	Firmicutes	Bacteria
504	<i>Bacillus pumilus</i> SAFR-032	Firmicutes	Bacteria
505	<i>Bacillus megaterium</i> QM B1551	Firmicutes	Bacteria
506	<i>Macrococcus caseolyticus</i> JCSC5402	Firmicutes	Bacteria
507	<i>Staphylococcus saprophyticus</i> ssp. <i>saprophyticus</i> ATCC 15305	Firmicutes	Bacteria
508	<i>Staphylococcus lugdunensis</i> HKU09-01	Firmicutes	Bacteria
509	<i>Staphylococcus haemolyticus</i> JCSC1435	Firmicutes	Bacteria
510	<i>Staphylococcus epidermidis</i> RP62A	Firmicutes	Bacteria
511	<i>Staphylococcus carnosus</i> ssp. <i>carnosus</i> TM300	Firmicutes	Bacteria
512	<i>Staphylococcus aureus</i> ssp. <i>aureus</i> NCTC 8325	Firmicutes	Bacteria
513	<i>Streptobacillus moniliformis</i> DSM 12112	Fusobacteria	Bacteria
514	<i>Sebaldella termitidis</i> ATCC 33386	Fusobacteria	Bacteria
515	<i>Leptotrichia buccalis</i> C-1013-b	Fusobacteria	Bacteria
516	<i>Fusobacterium nucleatum</i> ssp. <i>nucleatum</i> ATCC 25586	Fusobacteria	Bacteria
517	<i>Gemmatimonas aurantiaca</i> T-27	Gemmatimonadetes	Bacteria
518	<i>Thermodesulfovibrio yellowstonii</i> DSM 11347	Nitrospirae	Bacteria
519	<i>Rhodopirellula baltica</i> SH 1	Planctomycetes	Bacteria
520	<i>Pirellula staleyi</i> DSM 6068	Planctomycetes	Bacteria
521	<i>Nautilia profundicola</i> AmH	Proteobacteria	Bacteria
522	<i>Sulfurospirillum deleyianum</i> DSM 6946	Proteobacteria	Bacteria
523	<i>Arcobacter butzleri</i> RM4018	Proteobacteria	Bacteria
524	<i>Campylobacter hominis</i> ATCC BAA-381	Proteobacteria	Bacteria
525	<i>Campylobacter lari</i> RM2100	Proteobacteria	Bacteria
526	<i>Campylobacter curvus</i> 525.92	Proteobacteria	Bacteria
527	<i>Campylobacter concisus</i> 13826	Proteobacteria	Bacteria
528	<i>Campylobacter jejuni</i> ssp. <i>jejuni</i> NCTC 11168	Proteobacteria	Bacteria
529	<i>Campylobacter fetus</i> ssp. <i>fetus</i> 82-40	Proteobacteria	Bacteria
530	<i>Sulfurimonas denitrificans</i> DSM 1251	Proteobacteria	Bacteria
531	<i>Wolinella succinogenes</i> DSM 1740	Proteobacteria	Bacteria
532	<i>Helicobacter hepaticus</i> ATCC 51449	Proteobacteria	Bacteria
533	<i>Helicobacter mustelae</i> 12198	Proteobacteria	Bacteria

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
534	<i>Helicobacter acinonychis</i> Sheeba	Proteobacteria	Bacteria
535	<i>Helicobacter pylori</i> 26695	Proteobacteria	Bacteria
536	<i>Nitratiruptor</i> sp. SB155-2	Proteobacteria	Bacteria
537	<i>Sulfurovum</i> sp. NBC37-1	Proteobacteria	Bacteria
538	<i>Bdellovibrio bacteriovorus</i> HD100	Proteobacteria	Bacteria
539	<i>Syntrophus aciditrophicus</i> SB	Proteobacteria	Bacteria
540	<i>Syntrophobacter fumaroxidans</i> MPOB	Proteobacteria	Bacteria
541	<i>Desulfotalea psychrophila</i> LSv54	Proteobacteria	Bacteria
542	<i>Desulfatibacillum alkenivorans</i> AK-01	Proteobacteria	Bacteria
543	<i>Desulfobacterium autotrophicum</i> HRM2	Proteobacteria	Bacteria
544	<i>Desulfococcus oleovorans</i> Hxd3	Proteobacteria	Bacteria
545	<i>Desulfohalobium retbaense</i> DSM 5692	Proteobacteria	Bacteria
546	<i>Desulfomicrobium baculatum</i> DSM 4028	Proteobacteria	Bacteria
547	<i>Lawsonia intracellularis</i> PHE/MN1-00	Proteobacteria	Bacteria
548	<i>Desulfovibrio magneticus</i> RS-1	Proteobacteria	Bacteria
549	<i>Desulfovibrio vulgaris</i> Hildenborough	Proteobacteria	Bacteria
550	<i>Desulfovibrio salexigens</i> DSM 2638	Proteobacteria	Bacteria
551	<i>Desulfovibrio desulfuricans</i> ssp. <i>desulfuricans</i> G20	Proteobacteria	Bacteria
552	<i>Pelobacter propionicus</i> DSM 2379	Proteobacteria	Bacteria
553	<i>Pelobacter carbinolicus</i> DSM 2380	Proteobacteria	Bacteria
554	<i>Geobacter uraniireducens</i> Rf4	Proteobacteria	Bacteria
555	<i>Geobacter</i> sp. FRC-32	Proteobacteria	Bacteria
556	<i>Geobacter lovleyi</i> SZ	Proteobacteria	Bacteria
557	<i>Geobacter bemidjensis</i> Bem	Proteobacteria	Bacteria
558	<i>Geobacter sulfurreducens</i> PCA	Proteobacteria	Bacteria
559	<i>Geobacter metallireducens</i> GS-15	Proteobacteria	Bacteria
560	<i>Haliangium ochraceum</i> DSM 14365	Proteobacteria	Bacteria
561	<i>Sorangium cellulosum</i> So ce 56	Proteobacteria	Bacteria
562	<i>Anaeromyxobacter</i> sp. Fw109-5	Proteobacteria	Bacteria
563	<i>Anaeromyxobacter dehalogenans</i> 2CP-C	Proteobacteria	Bacteria
564	<i>Myxococcus xanthus</i> DK 1622	Proteobacteria	Bacteria
565	<i>Magnetococcus</i> sp. MC-1	Proteobacteria	Bacteria
566	<i>Sideroxydans lithotrophicus</i> ES-1	Proteobacteria	Bacteria
567	<i>Aromatoleum aromaticum</i> EbN1	Proteobacteria	Bacteria
568	<i>Dechloromonas aromatica</i> RCB	Proteobacteria	Bacteria
569	<i>Thauera</i> sp. MZ1T	Proteobacteria	Bacteria
570	<i>Laribacter hongkongensis</i> HLHK9	Proteobacteria	Bacteria
571	<i>Chromobacterium violaceum</i> ATCC 12472	Proteobacteria	Bacteria
572	<i>Neisseria meningitidis</i> Z2491	Proteobacteria	Bacteria
573	<i>Neisseria gonorrhoeae</i> FA 1090	Proteobacteria	Bacteria
574	<i>Methylothermobacter mobilis</i> JLW8	Proteobacteria	Bacteria
575	<i>Methylovorus</i> sp. SIP3-4	Proteobacteria	Bacteria
576	<i>Methylobacillus flagellatus</i> KT	Proteobacteria	Bacteria
577	<i>Thiobacillus denitrificans</i> ATCC 25259	Proteobacteria	Bacteria

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
578	<i>Candidatus Accumulibacter phosphatis</i> clade IIA UW-1	Proteobacteria	Bacteria
579	<i>Methylibium petroleiphilum</i> PM1	Proteobacteria	Bacteria
580	<i>Leptothrix cholodnii</i> SP-6	Proteobacteria	Bacteria
581	<i>Ralstonia eutropha</i> JMP134	Proteobacteria	Bacteria
582	<i>Cupriavidus taiwanensis</i>	Proteobacteria	Bacteria
583	<i>Cupriavidus metallidurans</i> CH34	Proteobacteria	Bacteria
584	<i>Ralstonia pickettii</i> 12J	Proteobacteria	Bacteria
585	<i>Ralstonia solanacearum</i> GMI1000	Proteobacteria	Bacteria
586	<i>Polynucleobacter necessarius</i> ssp. <i>asymbioticus</i> QLW-P1DMWA-1	Proteobacteria	Bacteria
587	<i>Burkholderia phytofirmans</i> PsJN	Proteobacteria	Bacteria
588	<i>Burkholderia phymatum</i> STM815	Proteobacteria	Bacteria
589	<i>Burkholderia thailandensis</i> E264	Proteobacteria	Bacteria
590	<i>Burkholderia pseudomallei</i> K96243	Proteobacteria	Bacteria
591	<i>Burkholderia mallei</i> ATCC 23344	Proteobacteria	Bacteria
592	<i>Burkholderia</i> sp. 383	Proteobacteria	Bacteria
593	<i>Burkholderia ambifaria</i> AMMD	Proteobacteria	Bacteria
594	<i>Burkholderia cenocepacia</i> AU 1054	Proteobacteria	Bacteria
595	<i>Burkholderia multivorans</i> ATCC 17616	Proteobacteria	Bacteria
596	<i>Burkholderia vietnamiensis</i> G4	Proteobacteria	Bacteria
597	<i>Burkholderia xenovorans</i> LB400	Proteobacteria	Bacteria
598	<i>Burkholderia glumae</i> BGR1	Proteobacteria	Bacteria
599	<i>Rhodoferax ferrireducens</i> T118	Proteobacteria	Bacteria
600	<i>Verminophrobacter eiseniae</i> EF01-2	Proteobacteria	Bacteria
601	<i>Delftia acidovorans</i> SPH-1	Proteobacteria	Bacteria
602	<i>Polaromonas</i> sp. JS666	Proteobacteria	Bacteria
603	<i>Polaromonas naphthalenivorans</i> CJ2	Proteobacteria	Bacteria
604	<i>Variovorax paradoxus</i> S110	Proteobacteria	Bacteria
605	<i>Acidovorax ebreus</i> TPSY	Proteobacteria	Bacteria
606	<i>Acidovorax</i> sp. JS42	Proteobacteria	Bacteria
607	<i>Acidovorax citrulli</i> AAC00-1	Proteobacteria	Bacteria
608	<i>Herminiimonas arsenicoxydans</i>	Proteobacteria	Bacteria
609	<i>Janthinobacterium</i> sp. Marseille	Proteobacteria	Bacteria
610	<i>Bordetella petrii</i> DSM 12804	Proteobacteria	Bacteria
611	<i>Bordetella avium</i> 197N	Proteobacteria	Bacteria
612	<i>Bordetella pertussis</i> Tohama I	Proteobacteria	Bacteria
613	<i>Bordetella parapertussis</i> 12822	Proteobacteria	Bacteria
614	<i>Bordetella bronchiseptica</i> RB50	Proteobacteria	Bacteria
615	<i>Nitrospira multififormis</i> ATCC 25196	Proteobacteria	Bacteria
616	<i>Nitrosomonas eutropha</i> C91	Proteobacteria	Bacteria
617	<i>Nitrosomonas europaea</i> ATCC 19718	Proteobacteria	Bacteria
618	<i>Caulobacter</i> sp. K31	Proteobacteria	Bacteria
619	<i>Caulobacter crescentus</i> CB15	Proteobacteria	Bacteria
620	<i>Caulobacter segnis</i> ATCC 21756	Proteobacteria	Bacteria
621	<i>Phenylobacterium zucineum</i> HLK1	Proteobacteria	Bacteria

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
622	<i>Erythrobacter litoralis</i> HTCC2594	Proteobacteria	Bacteria
623	<i>Sphingopyxis alaskensis</i> RB2256	Proteobacteria	Bacteria
624	<i>Novosphingobium aromaticivorans</i> DSM 12444	Proteobacteria	Bacteria
625	<i>Sphingobium japonicum</i> UT26S	Proteobacteria	Bacteria
626	<i>Sphingomonas wittichii</i> RW1	Proteobacteria	Bacteria
627	<i>Zymomonas mobilis</i> ssp. <i>mobilis</i> ZM4	Proteobacteria	Bacteria
628	<i>Maricaulis maris</i> MCS10	Proteobacteria	Bacteria
629	<i>Hirschia baltica</i> ATCC 49814	Proteobacteria	Bacteria
630	<i>Hyphomonas neptunium</i> ATCC 15444	Proteobacteria	Bacteria
631	<i>Dinoroseobacter shibae</i> DFL 12	Proteobacteria	Bacteria
632	<i>Jannaschia</i> sp. CCS1	Proteobacteria	Bacteria
633	<i>Ruegeria</i> sp. TM1040	Proteobacteria	Bacteria
634	<i>Ruegeria pomeroyi</i> DSS-3	Proteobacteria	Bacteria
635	<i>Roseobacter denitrificans</i> OCh 114	Proteobacteria	Bacteria
636	<i>Rhodobacter sphaeroides</i> 2.4.1	Proteobacteria	Bacteria
637	<i>Rhodobacter capsulatus</i> SB 1003	Proteobacteria	Bacteria
638	<i>Paracoccus denitrificans</i> PD1222	Proteobacteria	Bacteria
639	<i>Magnetospirillum magneticum</i> AMB-1	Proteobacteria	Bacteria
640	<i>Rhodospirillum centenum</i> SW	Proteobacteria	Bacteria
641	<i>Rhodospirillum rubrum</i> ATCC 11170	Proteobacteria	Bacteria
642	<i>Azospirillum</i> sp. B510	Proteobacteria	Bacteria
643	<i>Granulibacter bethesdensis</i> CGDNIH1	Proteobacteria	Bacteria
644	<i>Gluconacetobacter diazotrophicus</i> PAI 5	Proteobacteria	Bacteria
645	<i>Gluconobacter oxydans</i> 621H	Proteobacteria	Bacteria
646	<i>Acetobacter pasteurianus</i> IFO 3283-01	Proteobacteria	Bacteria
647	<i>Candidatus Puniceispirillum marinum</i> IMCC1322	Proteobacteria	Bacteria
648	<i>Candidatus Pelagibacter ubique</i> HTCC1062	Proteobacteria	Bacteria
649	<i>Neorickettsia sennetsu</i> Miyayama	Proteobacteria	Bacteria
650	<i>Neorickettsia risticii</i> Illinois	Proteobacteria	Bacteria
651	<i>Wolbachia</i> endosymbiont of <i>Culex quinquefasciatus</i> Pel	Proteobacteria	Bacteria
652	<i>Wolbachia</i> endosymbiont of <i>Drosophila melanogaster</i>	Proteobacteria	Bacteria
653	<i>Wolbachia</i> endosymbiont TRS of <i>Brugia malayi</i>	Proteobacteria	Bacteria
654	<i>Wolbachia</i> sp. wRi	Proteobacteria	Bacteria
655	<i>Ehrlichia chaffeensis</i> Arkansas	Proteobacteria	Bacteria
656	<i>Ehrlichia canis</i> Jake	Proteobacteria	Bacteria
657	<i>Ehrlichia ruminantium</i> Welgevonden	Proteobacteria	Bacteria
658	<i>Anaplasma phagocytophilum</i> HZ	Proteobacteria	Bacteria
659	<i>Anaplasma marginale</i> St. Maries	Proteobacteria	Bacteria
660	<i>Anaplasma centrale</i> Israel	Proteobacteria	Bacteria
661	<i>Orientia tsutsugamushi</i> Boryong	Proteobacteria	Bacteria
662	<i>Rickettsia bellii</i> RML369-C	Proteobacteria	Bacteria
663	<i>Rickettsia canadensis</i> McKiel	Proteobacteria	Bacteria
664	<i>Rickettsia typhi</i> Wilmington	Proteobacteria	Bacteria
665	<i>Rickettsia prowazekii</i> Madrid E	Proteobacteria	Bacteria

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
666	<i>Rickettsia peacockii</i> Rustic	Proteobacteria	Bacteria
667	<i>Rickettsia felis</i> URRWXCa2	Proteobacteria	Bacteria
668	<i>Rickettsia massiliae</i> MTU5	Proteobacteria	Bacteria
669	<i>Rickettsia africae</i> ESF-5	Proteobacteria	Bacteria
670	<i>Rickettsia akari</i> Hartford	Proteobacteria	Bacteria
671	<i>Rickettsia rickettsii</i> Sheila Smith	Proteobacteria	Bacteria
672	<i>Rickettsia conorii</i> Malish 7	Proteobacteria	Bacteria
673	<i>Xanthobacter autotrophicus</i> Py2	Proteobacteria	Bacteria
674	<i>Azorhizobium caulinodans</i> ORS 571	Proteobacteria	Bacteria
675	<i>Methylobacterium chloromethanicum</i> CM4	Proteobacteria	Bacteria
676	<i>Methylobacterium extorquens</i> PA1	Proteobacteria	Bacteria
677	<i>Methylobacterium</i> sp. 4-46	Proteobacteria	Bacteria
678	<i>Methylobacterium populi</i> BJ001	Proteobacteria	Bacteria
679	<i>Methylobacterium nodulans</i> ORS 2060	Proteobacteria	Bacteria
680	<i>Methylobacterium radiotolerans</i> JCM 2831	Proteobacteria	Bacteria
681	<i>Candidatus Hodgkinia cicadicola</i> Dsem	Proteobacteria	Bacteria
682	<i>Ochrobactrum anthropi</i> ATCC 49188	Proteobacteria	Bacteria
683	<i>Brucella microti</i> CCM 4915	Proteobacteria	Bacteria
684	<i>Brucella canis</i> ATCC 23365	Proteobacteria	Bacteria
685	<i>Brucella suis</i> 1330	Proteobacteria	Bacteria
686	<i>Brucella melitensis</i> bv. 1 16M	Proteobacteria	Bacteria
687	<i>Brucella ovis</i> ATCC 25840	Proteobacteria	Bacteria
688	<i>Brucella abortus</i> bv. 1 9-941	Proteobacteria	Bacteria
689	<i>Rhizobium</i> sp. NGR234	Proteobacteria	Bacteria
690	<i>Sinorhizobium medicae</i> WSM419	Proteobacteria	Bacteria
691	<i>Sinorhizobium meliloti</i> 1021	Proteobacteria	Bacteria
692	<i>Rhizobium etli</i> CFN 42	Proteobacteria	Bacteria
693	<i>Rhizobium leguminosarum</i> bv. <i>viciae</i> 3841	Proteobacteria	Bacteria
694	<i>Agrobacterium vitis</i> S4	Proteobacteria	Bacteria
695	<i>Agrobacterium radiobacter</i> K84	Proteobacteria	Bacteria
696	<i>Agrobacterium tumefaciens</i> C58	Proteobacteria	Bacteria
697	<i>Candidatus Liberibacter asiaticus</i> psy62	Proteobacteria	Bacteria
698	<i>Chelativorans</i> sp. BNC1	Proteobacteria	Bacteria
699	<i>Parvibaculum lavamentivorans</i> DS-1	Proteobacteria	Bacteria
700	<i>Mesorhizobium loti</i> MAFF303099	Proteobacteria	Bacteria
701	<i>Methylocella silvestris</i> BL2	Proteobacteria	Bacteria
702	<i>Beijerinckia indica</i> ssp. <i>indica</i> ATCC 9039	Proteobacteria	Bacteria
703	<i>Oligotropha carboxidovorans</i> OM5	Proteobacteria	Bacteria
704	<i>Rhodopseudomonas palustris</i> CGA009	Proteobacteria	Bacteria
705	<i>Nitrobacter winogradskyi</i> Nb-255	Proteobacteria	Bacteria
706	<i>Nitrobacter hamburgensis</i> XI4	Proteobacteria	Bacteria
707	<i>Bradyrhizobium</i> sp. ORS278	Proteobacteria	Bacteria
708	<i>Bradyrhizobium japonicum</i> USDA 110	Proteobacteria	Bacteria
709	<i>Bartonella tribocorum</i> CIP 105476	Proteobacteria	Bacteria

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
710	<i>Bartonella henselae</i> Houston-1	Proteobacteria	Bacteria
711	<i>Bartonella grahamii</i> as4aup	Proteobacteria	Bacteria
712	<i>Bartonella quintana</i> Toulouse	Proteobacteria	Bacteria
713	<i>Bartonella bacilliformis</i> KC583	Proteobacteria	Bacteria
714	<i>Acidithiobacillus ferrooxidans</i> ATCC 23270	Proteobacteria	Bacteria
715	<i>Mannheimia succiniciproducens</i> MBEL55E	Proteobacteria	Bacteria
716	<i>Aggregatibacter aphrophilus</i> NJ8700	Proteobacteria	Bacteria
717	<i>Aggregatibacter actinomycetemcomitans</i> D11S-1	Proteobacteria	Bacteria
718	<i>Haemophilus somnus</i> 129PT	Proteobacteria	Bacteria
719	<i>Pasteurella multocida</i> ssp. <i>multocida</i> Pm70	Proteobacteria	Bacteria
720	<i>Haemophilus parasuis</i> SH0165	Proteobacteria	Bacteria
721	<i>Haemophilus ducreyi</i> 35000HP	Proteobacteria	Bacteria
722	<i>Haemophilus influenzae</i> Rd KW20	Proteobacteria	Bacteria
723	<i>Actinobacillus succinogenes</i> 130Z	Proteobacteria	Bacteria
724	<i>Actinobacillus pleuropneumoniae</i> L20	Proteobacteria	Bacteria
725	<i>Tolumonas auensis</i> DSM 9187	Proteobacteria	Bacteria
726	<i>Aeromonas salmonicida</i> ssp. <i>salmonicida</i> A449	Proteobacteria	Bacteria
727	<i>Aeromonas hydrophila</i> ssp. <i>hydrophila</i> ATCC 7966	Proteobacteria	Bacteria
728	<i>Aliivibrio salmonicida</i> LFI1238	Proteobacteria	Bacteria
729	<i>Vibrio fischeri</i> ES114	Proteobacteria	Bacteria
730	<i>Vibrio parahaemolyticus</i> RIMD 2210633	Proteobacteria	Bacteria
731	<i>Vibrio harveyi</i> ATCC BAA-1116	Proteobacteria	Bacteria
732	<i>Vibrio</i> sp. Ex25	Proteobacteria	Bacteria
733	<i>Vibrio splendidus</i> LGP32	Proteobacteria	Bacteria
734	<i>Vibrio vulnificus</i> YJ016	Proteobacteria	Bacteria
735	<i>Vibrio cholerae</i> O1 biov. El Tor N16961	Proteobacteria	Bacteria
736	<i>Photobacterium profundum</i> SS9	Proteobacteria	Bacteria
737	<i>Psychromonas ingrahamii</i> 37	Proteobacteria	Bacteria
738	<i>Idiomarina loihiensis</i> L2TR	Proteobacteria	Bacteria
739	<i>Shewanella piezotolerans</i> WP3	Proteobacteria	Bacteria
740	<i>Shewanella loihica</i> PV-4	Proteobacteria	Bacteria
741	<i>Shewanella halifaxensis</i> HAW-EB4	Proteobacteria	Bacteria
742	<i>Shewanella sediminis</i> HAW-EB3	Proteobacteria	Bacteria
743	<i>Shewanella denitrificans</i> OS217	Proteobacteria	Bacteria
744	<i>Shewanella pealeana</i> ATCC 700345	Proteobacteria	Bacteria
745	<i>Shewanella oneidensis</i> MR-1	Proteobacteria	Bacteria
746	<i>Shewanella baltica</i> OS155	Proteobacteria	Bacteria
747	<i>Shewanella woodyi</i> ATCC 51908	Proteobacteria	Bacteria
748	<i>Shewanella</i> sp. MR-7	Proteobacteria	Bacteria
749	<i>Shewanella amazonensis</i> SB2B	Proteobacteria	Bacteria
750	<i>Shewanella violacea</i> DSS12	Proteobacteria	Bacteria
751	<i>Shewanella frigidimarina</i> NCIMB 400	Proteobacteria	Bacteria
752	<i>Shewanella putrefaciens</i> CN-32	Proteobacteria	Bacteria
753	<i>Colwellia psychrerythraea</i> 34H	Proteobacteria	Bacteria

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
754	<i>Pseudoalteromonas atlantica</i> T6c	Proteobacteria	Bacteria
755	<i>Pseudoalteromonas haloplanktis</i> TAC125	Proteobacteria	Bacteria
756	<i>Teredinibacter turnerae</i> T7901	Proteobacteria	Bacteria
757	<i>Saccharophagus degradans</i> 2-40	Proteobacteria	Bacteria
758	<i>Marinobacter aquaeolei</i> VT8	Proteobacteria	Bacteria
759	<i>Alteromonas macleodii</i> Deep ecotype	Proteobacteria	Bacteria
760	<i>Hahella chejuensis</i> KCTC 2396	Proteobacteria	Bacteria
761	<i>Kangiella koreensis</i> DSM 16069	Proteobacteria	Bacteria
762	<i>Alcanivorax borkumensis</i> SK2	Proteobacteria	Bacteria
763	<i>Marinomonas</i> sp. MWYL1	Proteobacteria	Bacteria
764	<i>Chromohalobacter salexigens</i> DSM 3043	Proteobacteria	Bacteria
765	<i>Methylococcus capsulatus</i> Bath	Proteobacteria	Bacteria
766	<i>Dichelobacter nodosus</i> VCS1703A	Proteobacteria	Bacteria
767	<i>Stenotrophomonas maltophilia</i> R551-3	Proteobacteria	Bacteria
768	<i>Xylella fastidiosa</i> 9a5c	Proteobacteria	Bacteria
769	<i>Xanthomonas axonopodis</i> pv. citri 306	Proteobacteria	Bacteria
770	<i>Xanthomonas albilineans</i>	Proteobacteria	Bacteria
771	<i>Xanthomonas oryzae</i> pv. oryzae KACC10331	Proteobacteria	Bacteria
772	<i>Xanthomonas campestris</i> pv. campestris ATCC 33913	Proteobacteria	Bacteria
773	<i>Halothiobacillus neapolitanus</i> c2	Proteobacteria	Bacteria
774	<i>Alkalilimnicola ehrlichii</i> MLHE-1	Proteobacteria	Bacteria
775	<i>Thioalkalivibrio</i> sp. HL-EbGR7	Proteobacteria	Bacteria
776	<i>Halorhodospira halophila</i> SL1	Proteobacteria	Bacteria
777	<i>Allochromatium vinosum</i> DSM 180	Proteobacteria	Bacteria
778	<i>Nitrosococcus halophilus</i> Nc4	Proteobacteria	Bacteria
779	<i>Nitrosococcus oceani</i> ATCC 19707	Proteobacteria	Bacteria
780	<i>Coxiella burnetii</i> RSA 493	Proteobacteria	Bacteria
781	<i>Legionella longbeachae</i> NSW150	Proteobacteria	Bacteria
782	<i>Legionella pneumophila</i> ssp. pneumophila Philadelphia 1	Proteobacteria	Bacteria
783	<i>Baumannia cicadellinicola</i> Hc	Proteobacteria	Bacteria
784	Candidatus <i>Carsonella ruddii</i> PV	Proteobacteria	Bacteria
785	Candidatus <i>Vesicomysocius okutanii</i> HA	Proteobacteria	Bacteria
786	Candidatus <i>Ruthia magnifica</i> Cm	Proteobacteria	Bacteria
787	<i>Cronobacter turicensis</i> z3032	Proteobacteria	Bacteria
788	<i>Cronobacter sakazakii</i> ATCC BAA-894	Proteobacteria	Bacteria
789	Candidatus <i>Riesia pediculicola</i> USDA	Proteobacteria	Bacteria
790	<i>Dickeya zea</i> Ech1591	Proteobacteria	Bacteria
791	<i>Dickeya dadantii</i> Ech703	Proteobacteria	Bacteria
792	Candidatus <i>Hamiltonella defensa</i> 5AT	Proteobacteria	Bacteria
793	Candidatus <i>Blochmannia floridanus</i>	Proteobacteria	Bacteria
794	<i>Pectobacterium wasabiae</i> WPP163	Proteobacteria	Bacteria
795	<i>Pectobacterium atrosepticum</i> SCRI1043	Proteobacteria	Bacteria
796	<i>Pectobacterium carotovorum</i> ssp. carotovorum PCI	Proteobacteria	Bacteria
797	<i>Sodalis glossinidius morsitans</i>	Proteobacteria	Bacteria

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
798	<i>Pantoea ananatis</i> LMG 20103	Proteobacteria	Bacteria
799	<i>Wigglesworthia glossinidia</i>	Proteobacteria	Bacteria
800	<i>Buchnera aphidicola</i> APS	Proteobacteria	Bacteria
801	<i>Photorhabdus asymbiotica</i>	Proteobacteria	Bacteria
802	<i>Photorhabdus luminescens</i> ssp. <i>laumondii</i> TTO1	Proteobacteria	Bacteria
803	<i>Edwardsiella ictaluri</i> 93-146	Proteobacteria	Bacteria
804	<i>Edwardsiella tarda</i> EIB202	Proteobacteria	Bacteria
805	<i>Yersinia pseudotuberculosis</i> IP 32953	Proteobacteria	Bacteria
806	<i>Yersinia pestis</i> CO92	Proteobacteria	Bacteria
807	<i>Yersinia enterocolitica</i> ssp. <i>enterocolitica</i> 8081	Proteobacteria	Bacteria
808	<i>Xenorhabdus bovienii</i> SS-2004	Proteobacteria	Bacteria
809	<i>Shigella sonnei</i> Ss046	Proteobacteria	Bacteria
810	<i>Shigella flexneri</i> 2a 2457T	Proteobacteria	Bacteria
811	<i>Shigella dysenteriae</i> Sd197	Proteobacteria	Bacteria
812	<i>Shigella boydii</i> Sb227	Proteobacteria	Bacteria
813	<i>Serratia proteamaculans</i> 568	Proteobacteria	Bacteria
814	<i>Salmonella enterica</i> ssp. <i>enterica</i> ser. <i>Typhimurium</i> LT2	Proteobacteria	Bacteria
815	<i>Proteus mirabilis</i> HI4320	Proteobacteria	Bacteria
816	<i>Klebsiella variicola</i> At-22	Proteobacteria	Bacteria
817	<i>Klebsiella pneumoniae</i> ssp. <i>pneumoniae</i> MGH 78578	Proteobacteria	Bacteria
818	<i>Escherichia fergusonii</i> ATCC 35469	Proteobacteria	Bacteria
819	<i>Escherichia coli</i> K-12 subMG1655	Proteobacteria	Bacteria
820	<i>Erwinia tasmaniensis</i> Et1/99	Proteobacteria	Bacteria
821	<i>Erwinia pyrifoliae</i> Ep1/96	Proteobacteria	Bacteria
822	<i>Erwinia amylovora</i> ATCC 49946	Proteobacteria	Bacteria
823	<i>Enterobacter</i> sp. 638	Proteobacteria	Bacteria
824	<i>Citrobacter rodentium</i> ICC168	Proteobacteria	Bacteria
825	<i>Citrobacter koseri</i> ATCC BAA-895	Proteobacteria	Bacteria
826	<i>Azotobacter vinelandii</i> DJ	Proteobacteria	Bacteria
827	<i>Pseudomonas entomophila</i> L48	Proteobacteria	Bacteria
828	<i>Pseudomonas syringae</i> pv. <i>tomato</i> DC3000	Proteobacteria	Bacteria
829	<i>Pseudomonas stutzeri</i> A1501	Proteobacteria	Bacteria
830	<i>Pseudomonas putida</i> KT2440	Proteobacteria	Bacteria
831	<i>Pseudomonas fluorescens</i> Pf-5	Proteobacteria	Bacteria
832	<i>Pseudomonas mendocina</i> ymp	Proteobacteria	Bacteria
833	<i>Pseudomonas aeruginosa</i> PAO1	Proteobacteria	Bacteria
834	<i>Cellvibrio japonicus</i> Ueda107	Proteobacteria	Bacteria
835	<i>Psychrobacter</i> sp. PRwf-1	Proteobacteria	Bacteria
836	<i>Psychrobacter arcticus</i> 273-4	Proteobacteria	Bacteria
837	<i>Psychrobacter cryohalolentis</i> K5	Proteobacteria	Bacteria
838	<i>Acinetobacter baumannii</i> ATCC 17978	Proteobacteria	Bacteria
839	<i>Acinetobacter</i> sp. ADP1	Proteobacteria	Bacteria
840	<i>Thiomicrospira crunogena</i> XCL-2	Proteobacteria	Bacteria
841	<i>Francisella philomiragia</i> ssp. <i>philomiragia</i> ATCC 25017	Proteobacteria	Bacteria

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
842	<i>Francisella tularensis</i> ssp. <i>tularensis</i> SCHU S4	Proteobacteria	Bacteria
843	<i>Brachyspira hyodysenteriae</i> WAI	Spirochaetes	Bacteria
844	<i>Leptospira borgpetersenii</i> ser. <i>Hardjo-bovis</i> L550	Spirochaetes	Bacteria
845	<i>Leptospira interrogans</i> ser. <i>Lai</i> 56601	Spirochaetes	Bacteria
846	<i>Leptospira biflexa</i> ser. <i>Patoc Patoc 1 (Paris)</i>	Spirochaetes	Bacteria
847	<i>Treponema pallidum</i> ssp. <i>pallidum</i> Nichols	Spirochaetes	Bacteria
848	<i>Treponema denticola</i> ATCC 35405	Spirochaetes	Bacteria
849	<i>Borrelia garinii</i> PBI	Spirochaetes	Bacteria
850	<i>Borrelia afzelii</i> PKo	Spirochaetes	Bacteria
851	<i>Borrelia burgdorferi</i> B31	Spirochaetes	Bacteria
852	<i>Borrelia recurrentis</i> AI	Spirochaetes	Bacteria
853	<i>Borrelia duttonii</i> Ly	Spirochaetes	Bacteria
854	<i>Borrelia turicatae</i> 91E135	Spirochaetes	Bacteria
855	<i>Borrelia hermsii</i> DAH	Spirochaetes	Bacteria
856	<i>Aminobacterium colombiense</i> DSM 12261	Synergistetes	Bacteria
857	<i>Thermanaerovibrio acidaminovorans</i> DSM 6589	Synergistetes	Bacteria
858	<i>Candidatus Phytoplasma mali</i>	Tenericutes	Bacteria
859	<i>Aster yellows witches-broom phytoplasma</i> AYWB	Tenericutes	Bacteria
860	<i>Onion yellows phytoplasma</i> OY-M	Tenericutes	Bacteria
861	<i>Acholeplasma laidlawii</i> PG-8A	Tenericutes	Bacteria
862	<i>Mesoplasma florum</i> LI	Tenericutes	Bacteria
863	<i>Ureaplasma parvum</i> ser. 3 ATCC 700970	Tenericutes	Bacteria
864	<i>Ureaplasma urealyticum</i> ser. 10 ATCC 33699	Tenericutes	Bacteria
865	<i>Mycoplasma mycoides</i> ssp. <i>mycoides</i> SC PG1	Tenericutes	Bacteria
866	<i>Mycoplasma capricolum</i> ssp. <i>capricolum</i> ATCC 27343	Tenericutes	Bacteria
867	<i>Mycoplasma crocodyli</i> MP145	Tenericutes	Bacteria
868	<i>Mycoplasma conjunctivae</i> HRC/581	Tenericutes	Bacteria
869	<i>Mycoplasma penetrans</i> HF-2	Tenericutes	Bacteria
870	<i>Mycoplasma mobile</i> 163K	Tenericutes	Bacteria
871	<i>Mycoplasma arthritidis</i> 158L3-1	Tenericutes	Bacteria
872	<i>Mycoplasma agalactiae</i> PG2	Tenericutes	Bacteria
873	<i>Mycoplasma synoviae</i> 53	Tenericutes	Bacteria
874	<i>Mycoplasma pulmonis</i> UAB CTIP	Tenericutes	Bacteria
875	<i>Mycoplasma pneumoniae</i> M129	Tenericutes	Bacteria
876	<i>Mycoplasma hyopneumoniae</i> 232	Tenericutes	Bacteria
877	<i>Mycoplasma hominis</i>	Tenericutes	Bacteria
878	<i>Mycoplasma genitalium</i> G37	Tenericutes	Bacteria
879	<i>Mycoplasma gallisepticum</i> R(low)	Tenericutes	Bacteria
880	<i>Kosmotoga olearia</i> TBF 19.5.1	Thermotogae	Bacteria
881	<i>Petrotoga mobilis</i> SJ95	Thermotogae	Bacteria
882	<i>Fervidobacterium nodosum</i> Rt17-B1	Thermotogae	Bacteria
883	<i>Thermosipho melanesiensis</i> BI429	Thermotogae	Bacteria
884	<i>Thermosipho africanus</i> TCF52B	Thermotogae	Bacteria
885	<i>Thermotoga lettingae</i> TMO	Thermotogae	Bacteria

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
886	<i>Thermotoga sp. RQ2</i>	Thermotogae	Bacteria
887	<i>Thermotoga naphthophila</i> RKU-10	Thermotogae	Bacteria
888	<i>Thermotoga petrophila</i> RKU-1	Thermotogae	Bacteria
889	<i>Thermotoga neapolitana</i> DSM 4359	Thermotogae	Bacteria
890	<i>Thermotoga maritima</i> MSB8	Thermotogae	Bacteria
891	<i>Coralimargarita akajimensis</i> DSM 45221	Verrucomicrobia	Bacteria
892	<i>Opitutus terrae</i> PB90-1	Verrucomicrobia	Bacteria
893	<i>Methylacidiphilum infernorum</i> V4	Verrucomicrobia	Bacteria
894	<i>Akkermansia muciniphila</i> ATCC BAA-835	Verrucomicrobia	Bacteria
895	<i>Thermobaculum terrenum</i> ATCC BAA-798		Bacteria
896	<i>Hyperthermus butylicus</i> DSM 5456	Crenarchaeota	Archaea
897	<i>Aeropyrum pernix</i> K1	Crenarchaeota	Archaea
898	<i>Ignicoccus hospitalis</i> KIN4/I	Crenarchaeota	Archaea
899	<i>Staphylothermus marinus</i> F1	Crenarchaeota	Archaea
900	<i>Desulfurococcus kamchatkensis</i> 1221n	Crenarchaeota	Archaea
901	<i>Metallosphaera sedula</i> DSM 5348	Crenarchaeota	Archaea
902	<i>Sulfolobus tokodaii</i> 7	Crenarchaeota	Archaea
903	<i>Sulfolobus islandicus</i> Y.N.15.51	Crenarchaeota	Archaea
904	<i>Sulfolobus solfataricus</i> P2	Crenarchaeota	Archaea
905	<i>Sulfolobus acidocaldarius</i> DSM 639	Crenarchaeota	Archaea
906	<i>Thermofilum pendens</i> Hrk 5	Crenarchaeota	Archaea
907	<i>Caldivirga maquilingensis</i> IC-167	Crenarchaeota	Archaea
908	<i>Pyrobaculum calidifontis</i> JCM 11548	Crenarchaeota	Archaea
909	<i>Pyrobaculum arsenaticum</i> DSM 13514	Crenarchaeota	Archaea
910	<i>Pyrobaculum aerophilum</i> IM2	Crenarchaeota	Archaea
911	<i>Pyrobaculum islandicum</i> DSM 4184	Crenarchaeota	Archaea
912	<i>Thermoproteus neutrophilus</i> V24Sta	Crenarchaeota	Archaea
913	<i>Methanocella paludicola</i> SANAE	Euryarchaeota	Archaea
914	<i>Methanosaeta thermophila</i> PT	Euryarchaeota	Archaea
915	<i>Methanococcoides burtonii</i> DSM 6242	Euryarchaeota	Archaea
916	<i>Methanosarcina acetivorans</i> C2A	Euryarchaeota	Archaea
917	<i>Methanosarcina mazei</i> Go1	Euryarchaeota	Archaea
918	<i>Methanosarcina barkeri</i> Fusaro	Euryarchaeota	Archaea
919	<i>Methanohalophilus mahii</i> DSM 5219	Euryarchaeota	Archaea
920	<i>Methanosphaerula palustris</i> E1-9c	Euryarchaeota	Archaea
921	<i>Candidatus Methanoregula boonei</i> 6A8	Euryarchaeota	Archaea
922	<i>Methanospirillum hungatei</i> JF-1	Euryarchaeota	Archaea
923	<i>Methanocorpusculum labreanum</i> Z	Euryarchaeota	Archaea
924	<i>Methanoculleus marisnigri</i> JR1	Euryarchaeota	Archaea
925	<i>Methanopyrus kandleri</i> AV19	Euryarchaeota	Archaea
926	<i>Ferroglobus placidus</i> DSM 10642	Euryarchaeota	Archaea
927	<i>Archaeoglobus profundus</i> DSM 5631	Euryarchaeota	Archaea
928	<i>Archaeoglobus fulgidus</i> DSM 4304	Euryarchaeota	Archaea
929	<i>Thermococcus onnurineus</i> NA1	Euryarchaeota	Archaea

Table S5. Cont.

No.	Genome Name	Phyla/Kingdom	Superkingdom
930	<i>Thermococcus kodakarensis</i> KOD1	Euryarchaeota	Archaea
931	<i>Thermococcus gammatolerans</i> EJ3	Euryarchaeota	Archaea
932	<i>Thermococcus sibiricus</i> MM 739	Euryarchaeota	Archaea
933	<i>Pyrococcus horikoshii</i> OT3	Euryarchaeota	Archaea
934	<i>Pyrococcus abyssi</i> GE5	Euryarchaeota	Archaea
935	<i>Pyrococcus furiosus</i> DSM 3638	Euryarchaeota	Archaea
936	<i>Thermoplasma volcanium</i> GSS1	Euryarchaeota	Archaea
937	<i>Thermoplasma acidophilum</i> DSM 1728	Euryarchaeota	Archaea
938	<i>Picrophilus torridus</i> DSM 9790	Euryarchaeota	Archaea
939	<i>Haloquadratum walsbyi</i> DSM 16790	Euryarchaeota	Archaea
940	<i>Halomicrobium mukohataei</i> DSM 12286	Euryarchaeota	Archaea
941	<i>Halorhabdus utahensis</i> DSM 12940	Euryarchaeota	Archaea
942	<i>Haloterrigena turkmenica</i> DSM 5511	Euryarchaeota	Archaea
943	<i>Natronomonas pharaonis</i> DSM 2160	Euryarchaeota	Archaea
944	<i>Natrialba magadii</i> ATCC 43099	Euryarchaeota	Archaea
945	<i>Halorubrum lacusprofundi</i> ATCC 49239	Euryarchaeota	Archaea
946	<i>Haloferax volcanii</i> DS2	Euryarchaeota	Archaea
947	<i>Halobacterium salinarum</i> R1	Euryarchaeota	Archaea
948	<i>Halobacterium</i> sp. NRC-1	Euryarchaeota	Archaea
949	<i>Haloarcula marismortui</i> ATCC 43049	Euryarchaeota	Archaea
950	<i>Methanocaldococcus</i> sp. FS406-22	Euryarchaeota	Archaea
951	<i>Methanocaldococcus fervens</i> AG86	Euryarchaeota	Archaea
952	<i>Methanocaldococcus vulcanius</i> M7	Euryarchaeota	Archaea
953	<i>Methanocaldococcus jannaschii</i> DSM 2661	Euryarchaeota	Archaea
954	<i>Methanococcus aeolicus</i> Nankai-3	Euryarchaeota	Archaea
955	<i>Methanococcus maripaludis</i> S2	Euryarchaeota	Archaea
956	<i>Methanococcus vannieli</i> SB	Euryarchaeota	Archaea
957	<i>Methanothermobacter thermautotrophicus</i> Delta H	Euryarchaeota	Archaea
958	<i>Methanosphaera stadtmanae</i> DSM 3091	Euryarchaeota	Archaea
959	<i>Methanobrevibacter ruminantium</i> M1	Euryarchaeota	Archaea
960	<i>Methanobrevibacter smithii</i> ATCC 35061	Euryarchaeota	Archaea
961	uncultured methanogenic archaeon RC-1	Euryarchaeota	Archaea
962	<i>Aciduliprofundum boonei</i> T469	Euryarchaeota	Archaea
963	<i>Candidatus Korarchaeum cryptofilum</i> OPF8	Korarchaeota	Archaea
964	<i>Nanoarchaeum equitans</i> Kin4-M	Nanoarchaeota	Archaea
965	<i>Nitrosopumilus maritimus</i> SCM1	Thaumarchaeota	Archaea