

1 ***fusionDB*: assessing microbial diversity and**
2 **environmental preferences via functional similarity**
3 **networks**

4 **Zhu C^{1*}, Mahlich Y^{2*}, Miller M^{2*}, Bromberg Y^{1,2}**

5 ¹Department of Biochemistry and Microbiology, Rutgers University, 76 Lipman Dr, New
6 Brunswick, NJ 08873, USA

7 ²Institute for Advanced Study, Technische Universität München, Lichtenbergstrasse 2 a, D-
8 85748 Garching, Germany.

9
10 * Corresponding authors: czhu@bromberglab.org, ymahlich@bromberglab.org

11 Tel: +1.848.932.5638; Fax +1.848.932.8965

12 + These authors contributed equally to this manuscript.

13
14
15
16
17
18

Abstract

1
2 Microbial functional diversification is driven by environmental factors, *i.e.*
3 microorganisms inhabiting the same environmental niche tend to be more
4 functionally similar than those from different environments. In some cases, even
5 closely phylogenetically related microbes differ more across environments than
6 across taxa. While microbial similarities are often reported in terms of taxonomic
7 relationships, no existing databases directly links microbial functions to the
8 environment. We previously developed a method for comparing microbial functional
9 similarities on the basis of proteins translated from the sequenced genomes. Here we
10 describe *fusionDB*, a novel database that uses our functional data to represent 1,374
11 taxonomically distinct bacteria annotated with available metadata: habitat/niche,
12 preferred temperature, and oxygen use. Each microbe is encoded as a set of
13 functions represented by its proteome and individual microbes are connected via
14 common functions. Users can search *fusionDB* via combinations of organism names
15 and metadata. Moreover, the web interface allows mapping new microbial genomes
16 to the functional spectrum of reference bacteria, rendering interactive similarity
17 networks that highlight shared functionality. *fusionDB* provides a fast means of
18 comparing microbes, identifying potential horizontal gene transfer events, and
19 highlighting key environment-specific functionality.

20 *fusionDB* is publicly available at <http://services.bromberglab.org/fusiondb/>.

21

Introduction

1
2 Microorganisms are capable of carrying out much of molecular functionality relevant
3 to a range of human interests, including health, industrial production, and
4 bioremediation. Experimental study of these microbes to optimize their uses is
5 expensive and time-consuming; e.g. as many as three hundred
6 biochemical/physiological tests only reflect 5-20% of the bacterial functional potential
7 (Garrity GM, 2001). The recent drastic increase in the number of sequenced
8 microbial genomes has facilitated access to microbial molecular functionality from the
9 gene/protein sequence side, via databases like Pfam (Sayers, et al., 2009), COG
10 (Tatusov, et al., 2003), TIGRfam (Haft, et al., 2003), RAST (Aziz, et al., 2008) and
11 others. Note that the relatively low number of available experimental functional
12 annotations limits the power of these databases in recognizing microbial proteins that
13 provide novel functionality. Additional information about microbial environmental
14 preferences can be found, e.g. in GOLD (Pagani, et al., 2012). While it is well known
15 that environmental factors play an important role in microbial functionality (Cohan,
16 2001), none of the existing resources directly link environmental data to microbial
17 function.

18 We mapped bacterial proteins to molecular functions and studied the functional
19 relationships between bacteria in the light of their chosen habitats. We previously
20 developed *fusion* (Zhu, et al., 2015), an organism functional similarity network, which
21 can be used to broadly summarize the environmental factors driving microbial
22 functional diversification. Here we describe *fusionDB* – a database relating bacterial
23 *fusion* functional repertoires to the corresponding environmental niches. *fusionDB* is
24 explorable via a web-interface by querying for combinations of organism names and
25 environments. Users can also map new organism proteomes to the functional
26 repertoires of the reference organisms in *fusionDB*; including, notably, matching
27 proteins of yet unannotated function across organisms. The submitted organisms are
28 visualized, and can be further explored, interactively as *fusion* networks in the
29 context of selected reference genomes. Additionally, the web interface generates
30 *fusion+* networks, i.e. views that explicitly indicate shared microbial functions.

1 Our overall analyses of the *fusionDB* data for the first time give quantitative support
2 for the fact that environmental factors driving microbial functional diversification. To
3 demonstrate *fusionDB* functionality, for individual organisms we mapped a recently
4 sequenced genome of a freshwater *Synechococcus* bacterium to *fusionDB*. In line
5 with our previous findings (Zhu, et al., 2015), we demonstrate that this microorganism
6 is more functionally related to other fresh water Cyanobacteria than to the marine
7 *Synechococcus*. In a case study on *Bacillus* microbes we use *fusionDB* to track
8 organism-unique functions and illustrate the detection of core-function repertoires
9 that capture traces of environmentally driven horizontal gene transfer (HGT).
10 *fusionDB* is a unique tool that provides an easy way of analysing the, often
11 unannotated, molecular function spectrum of a given microbe. It further places this
12 microbe into a context of other reference organisms and relates the identified
13 microbial function to the preferred environmental conditions. Our approach allows for
14 detection of microbial functional similarities, often mediated via horizontal gene
15 transfer, that are difficult to recover via phylogenetic analysis. We note that *fusionDB*
16 may also be useful for the analysis of functional potentials encoded in microbiome
17 metagenomes. We expect that *fusionDB* will facilitate the study of environment-
18 specific microbial molecular functionalities, leading to improved understanding of
19 microbial lifestyles and to an increased number of applied bacterial uses.

Methods

Database setup. *fusionDB* is based on alignments of 4,284,540 proteins from 1,374 bacterial genomes (Dec. 2011 NCBI GenBank (Benson, et al., 2009). For each bacterium, we store its (1) NCBI taxonomic information (Sayers, et al., 2009) and, where available, (2) environmental metadata (temperature, oxygen requirements, and habitat; GOLD (Pagani, et al., 2012). The environments are generalized, e.g. *thermophiles* include hyper-thermophiles. “No data” is used to indicate missing annotations (Supplementary Online Material, SOM_Table 1). The general *fusion* (functional repertoire similarity-based organism network) protocol is described in (Zhu, et al., 2015). Briefly, all proteins in our database are aligned against each other using three iterations of PSI-BLAST (Altschul, et al., 1990) and the alignment length and sequence identity are used to compute HSSP (Rost, 2002). A network of protein similarities is then clustered using MCL (Dongen, 2000) clustering. For *fusionDB* the original *fusion* algorithm was modified to use less stringent protein functional similarity criteria (with HSSP distance cutoff = 10), which resulted in 457,576 functions (protein clusters; SOM Table 2). Each bacterium was thus mapped to a set of functions, its functional repertoire. Therefore our functional repertoires include all the bacterial functions, regardless of annotation. We are thus able to make function predictions, even the functions that have not been annotated before, for proteins in new bacteria.

Web interface. *fusionDB* web interface has two functions: *explore* and *map new organisms*. The *explore* section contains access to all the 1,374 bacteria and their metadata. Users can search these with (combinations of) organism names and environmental preferences by using text box input or built in filters. User-selected organism set is then used to create a *fusion* network, in which organism nodes are connected by functional similarity edges. The *fusion* network can be viewed in an interactive display, as well as downloaded as network data files or static images. The user-defined color labels of the organism nodes reflect microbial taxonomy or environment. In the interactive display clicking an organism node reveals its taxonomic information and environmental preferences, while clicking an edge

1 between two organisms yields a list of their shared functions. A *fusion+* network can
2 further be generated from the same list of organisms. There are two types of vertices
3 (nodes) in *fusion+*: organism nodes and function nodes. Organism nodes are
4 connected to each other only through the function nodes they share. The number of
5 edges (degree) of an organism node represents the total number of functions of the
6 organism; the relative position of each organism node is determined by the pull
7 *towards* other organisms via the common functions and *away* from others via unique
8 functions (Zhu, et al., 2015). Like *fusion*, *fusion+* can be interactively displayed,
9 downloaded, and colored by the users' choices. For both network types, users can
10 further retrieve the functions shared by the selected organisms - the core-functional
11 repertoire of the set. Note that the function annotation is from myRAST (including
12 hypothetical function or unknown function; Aziz, et al., 2008). This feature is an
13 efficient tool for investigating functions underlying organism diversification,
14 particularly within different environment conditions.

15 In the *map* section, users can submit their own new organism proteomes (in fasta
16 format) to our server. The submitted proteins are PSI-BLASTed against *fusionDB*
17 and assigned to stored functions using the HSSP distance cutoff = 10. Note that
18 novel proteins that can't be assigned to existing functional groups are reported as
19 functional singletons. Additionally, protein alignments that exceed 12 CPU hours of
20 run-time are eliminated from future consideration. In testing, we found that no more
21 than 0.1% of the proteins fall into this category. Although long run-times usually
22 indicate that query proteins likely align to many others in our database, they
23 contribute only a small fraction to the overall bacterial similarity and are eliminated for
24 the sake of a faster result turn-around. The server sends out emails to users when
25 mapping is finished. The *map* result page contains two tables: one is the list of
26 functions of the submitted bacterium, while the other contains pairwise functional
27 similarities (Eqn. 1) between the submitted bacterium and the reference proteomes in
28 *fusionDB* (SOM Figure 1).

29
$$\text{similarity} = \frac{\text{shared functions}}{\text{the larger functional repertoire size}} \quad \text{Eqn. 1}$$

1 Both tables can be easily sorted, searched and exported as comma-separated files.
2 The submitted proteome is further mapped to user-selected reference organisms with
3 *fusion* and/or *fusion+* as describe above.

4 ***Analysis of environment-driven organism similarity.*** For each environmental
5 condition in *fusionDB*, we sampled organism pairs where organisms were from (1)
6 the same condition (SC, e.g. both mesophiles) and (2) different conditions (DC, e.g.
7 thermophile vs. mesophile). To alleviate the effects of data bias, the organisms in
8 one pair were always selected from different taxonomic groups (different Families).
9 The smallest available set of pairs, SC-psychrophile contained 33 organisms from 17
10 Families (SOM_Table 1; 136 pairs – 48 same phylum, 88 different phyla; due to high
11 functional diversity of *Proteobacteria*, its classes were considered independent phyla).
12 For all other environment factors we sampled, 100 bootstrap times, 136 organism
13 pairs for both SC and DC sets, covering this same minimum taxonomic diversity. We
14 calculated the pairwise functional similarity (Eqn. 1) distributions and discarded
15 organism pairs with less than 5% similarity.

16

Results and Discussion

Map new *Synechococcus* genomes to *fusionDB*. We downloaded the full genome of *Synechococcus* sp. PCC 7502 (GCA_000317085.1) as translated protein sequence fasta (.faa file) from the NCBI Genbank (Benson, et al., 2009) and submitted it to our web interface. This 3,318 protein fresh water Cyanobacteria is isolated from a Sphagnum (peat moss) bog (Pagani, et al., 2012). 2,889 (87%) of the bacterial proteins mapped to 2,206 *fusionDB* functions and 426 (13%) were functional singletons; three proteins exceeded runtime and were excluded, (Methods). The whole process from submission to receiving a results notification e-mail took a little under three and a half hours. The mapping indicates that *Synechococcus* sp. PCC 7502 is functionally most similar (56%) to *Synechocystis* PCC 6803, a fresh water organism closely related to *Synechococcus*. It also shares a high functional similarity with a mud *Synechococcus* (*S.sp.* PCC 7002; 53%) and with other fresh water *Synechococcus* (*S.elongatus* PCC 7942 and *S.elongatus* PCC 6301; 52%). Notably, but not surprisingly, *Synechococcus* sp. PCC 7502 shares much less functional similarity (40-42%) with the marine *Synechococcus* bacteria. This relationship is clearly demonstrated by the *fusion+* networks (Figure 1). There are 874 functions shared by all the twelve *Synechococcus* (SOM Table 3), the core-function repertoire for this genus, and 1,128 functions shared among the fresh water *Synechococcus* (SOM Table 4). These additional 254 functions (SOM Table 5) are likely important for surviving in the fresh water, as opposed to the marine, environment, e.g. low salinity and low osmotic pressure.

Environment significantly affects microbial function. Not surprisingly, the SC-thermophile and SC-psychrophile pairs demonstrate significantly higher similarities comparing to all DC pairs (Figure 2A). Notably, the higher functional similarity between thermophiles than between psychrophiles suggests that protein functional adaptation to low temperature is less drastic than to high temperature – an interesting finding itself. Contrast to the extremophiles, mesophile organisms seem to have huge functional diversity as the SC-mesophile similarities are comparable to those the DC pairs (Figure 2A).

1 Different molecular pathways of aerobic-respiration and anaerobic-
2 respiration/fermentation explain the highest dissimilarity between the aerobes and
3 anaerobes (DC-anaerobe-aerobe; Figure 2B). Interestingly, the SC-anaerobe
4 similarities are higher than the SC-aerobe similarities, probably because the more
5 ancient anaerobic-respiration/fermentation machinery is more simple and conserved.

6 Habitat-based DC samples show lower pairwise organism similarity than SC samples
7 as well (Figure 2C), except for DC-fresh water-marine, which is not surprising due the
8 same aquatic condition. SC-host displays the lowest mean organism similarity of the
9 habitat SC samples. We speculate it is the result from the evolutionary pressure to
10 deal with diverse host defence mechanisms (Hornef, et al., 2002). The soil organisms
11 also share low functional similarity, which is likely due to soil's heterogeneity at
12 physical, chemical, and biological levels, from nano- to landscape scale (Bastian, et
13 al., 2009).

14 In general, SC organisms across all environmental factors are more functionally
15 similar than DC organisms (Figure 2; with exceptions mentioned above; Kolmogorov-
16 Smirnov test $p\text{-val} < 2.5e\text{-}6$). In other words, organisms in the same environment are
17 generally more similar than organisms from different environments. This finding is
18 intuitive and many studies have shown HGT within environment-specific microbiomes
19 (Kim, et al., 2012; Liu, et al., 2012; Saye, et al., 1987). Our results, however, for the
20 first time quantify on a broad scale the environmental impact on microorganism
21 function diversification.

22 **Case study of a temperature driven HGT event.** In *fusionDB explore*, we extracted
23 thermophilic, mesophilic, and psychrophilic species representatives (one per species)
24 of the *Bacillus* genus from *fusionDB*. We also added two other thermophilic
25 organisms, *D. carboxydivorans* CO-1-SRB and *S. acidophilus* TPY, to generate a
26 *fusion+* network (SOM_Table 6, SOM_Figure 2). The non-*Bacillus* thermophiles were
27 more closely related to the thermophilic *Bacilli*. All five thermophiles exclusively share
28 three functions. One is a likely pyruvate phosphate dikinase (PPDK) that, in
29 extremophiles, works as a primary glycolysis enzyme (Chastain, et al., 2011).
30 Phylogenetic analysis (SOM Methods) suggests an HGT event between thermophilic

1 organisms or a differential gene-loss in Bacilli that no longer live under high
2 temperature (SOM Figure 3). The other two shared functions are carried out by
3 proteins translated from mobile genetic elements (MGEs) that mediate the movement
4 of DNA within genomes or between bacteria (Frost, et al., 2005). Shared closely
5 related MGEs in distant organisms imply HGT (Krupovic, et al., 2013). We thus
6 suggest that *fusionDB* offers a fast and easy way to trace functionally-necessary
7 HGT within niche-specific microbial communities.

8 We have highlighted the importance of environmental factors for microbial function,
9 and demonstrated the capability of *fusionDB* to not only annotate functions, but also
10 directly link function to environment. Although it was developed for mapping new
11 microbial genomes, *fusionDB* also has the potential for microbiome annotation. By
12 mapping the proteins translated from metagenomes assembly to *fusionDB*, both the
13 functional and taxonomical can be obtained. We look forward to making *fusionDB*
14 more useful in this direction.

15

4. Conclusions

fusionDB links microbial functional similarities and environmental preferences. Our data analysis reveals environmental factors driving microbial functional diversification. Mapping new genomes to the reference genomes, it offers a novel, fast, and simple way to detect core-function repertoires, unique functions, as well as traces of HGT. With more microbial genome sequencing and further manual curation of environmental metadata, we expect that *fusionDB* will become an integral part of microbial functional analysis protocols in the near future.

Acknowledgements

We thank Drs. Burkhard Rost (TU Munich), Max Haggblom and Tamar Barkay (both Rutgers), and Tom O. Delmont (U Chicago) for all discussions and to those who deposit their data in public databases.

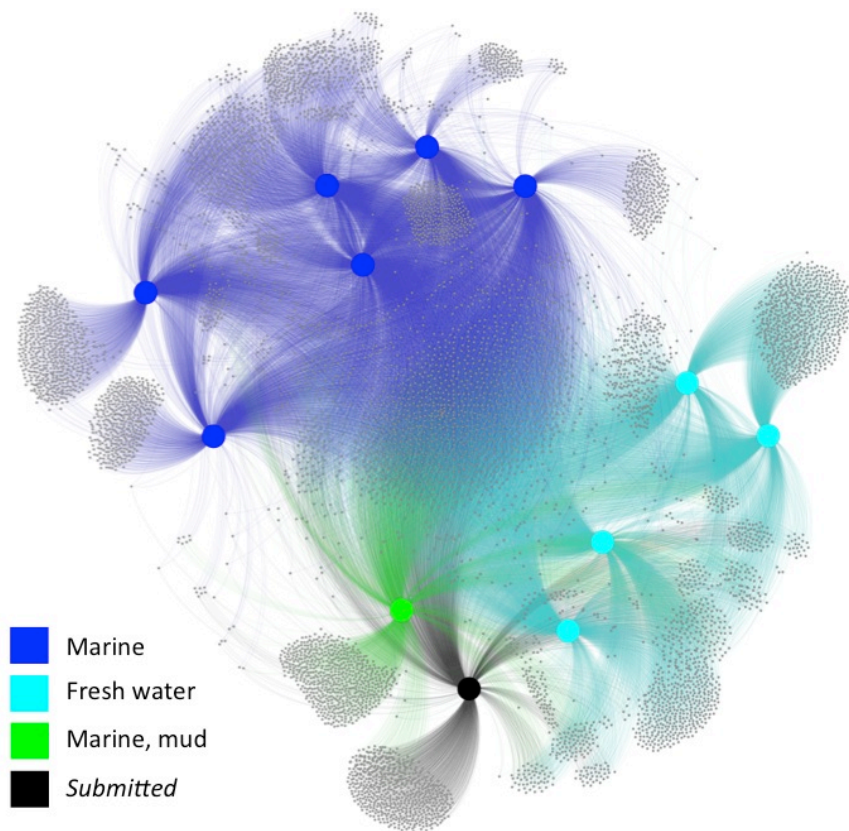
Funding

This work was supported by the NSF CAREER Award 1553289, USDA-NIFA 1015:0228906 and the TU München – Institute for Advanced Study Hans Fischer Fellowship, funded by the German Excellence Initiative and the EU Seventh Framework Programme, grant agreement 291763.

Conflict of Interest: none declared

1

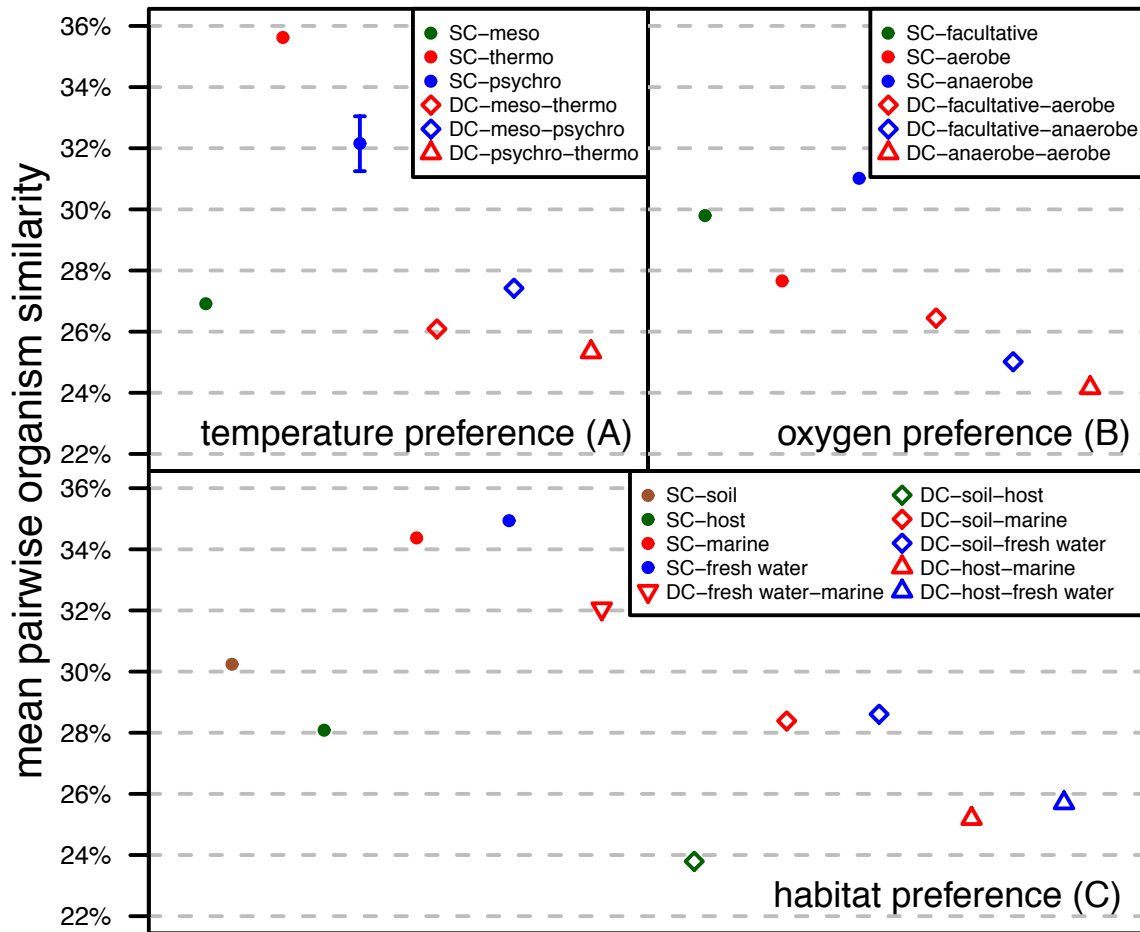
Figures



2

3 **Figure 1. The fusion+ view of all *Synechococcus* genomes.** The submitted
4 *Synechococcus sp.* PCC 7502 (black) cluster with the fresh water *Synechococcus*
5 organisms (light blue). Note that the *Synechococcus sp.* PCC 7002 (green), which is
6 isolated from marine mud, is salt tolerant but does not require salt for growth (see
7 (Zhu, et al., 2015)).

8



1
2 **Figure 2. Organism pairwise similarity is higher among organisms living in the**
3 **same environmental conditions.** The mean pairwise similarity for same (SC) and
4 different (DC) condition organisms according to (A) temperature preference, (B)
5 oxygen requirement, and (C) habitat. For all points without error bars, the standard
6 errors are vanishingly small.
7

References

- 1
- 2 Altschul, S.F., *et al.* (1990) Basic local alignment search tool, *J Mol Biol*, **215**, 403-410.
- 3 Aziz, R.K., *et al.* (2008) The RAST Server: rapid annotations using subsystems technology,
4 *BMC Genomics*, **9**, 75.
- 5 Bastian, M., Heymann, S. and Jacomy, M. (2009) *Gephi: An Open Source Software for*
6 *Exploring and Manipulating Networks*. 2009.
- 7 Benson, D.A., *et al.* (2009) GenBank, *Nucleic Acids Res*, **37**, D26-31.
- 8 Chastain, C.J., *et al.* (2011) Functional evolution of C4 pyruvate,orthophosphate dikinase,
9 *Journal of Experimental Botany*.
- 10 Cohan, F.M. (2001) Bacterial Species and Speciation, *Systematic Biology*, **50**, 513-524.
- 11 Dongen, S.v. (2000) Graph Clustering by Flow Simulation., *PhD thesis, University of Utrecht*.
- 12 Frost, L.S., *et al.* (2005) Mobile genetic elements: the agents of open source evolution, *Nat*
13 *Rev Micro*, **3**, 722-732.
- 14 Garrity GM, B.D., Castenholz RW, editors (2001) *Bergey's Manual of Systematic*
15 *Bacteriology, Volume 1*. Springer, New York (NY).
- 16 Haft, D.H., Selengut, J.D. and White, O. (2003) The TIGRFAMs database of protein families,
17 *Nucleic Acids Res*, **31**, 371-373.
- 18 Hornef, M.W., *et al.* (2002) Bacterial strategies for overcoming host innate and adaptive
19 immune responses, *Nat Immunol*, **3**, 1033-1040.
- 20 Kim, S.E., *et al.* (2012) Monitoring of horizontal gene transfer from agricultural
21 microorganisms to soil bacteria and analysis of microbial community in soils, *Journal of*
22 *microbiology and biotechnology*, **22**, 563-566.
- 23 Krupovic, M., *et al.* (2013) Insights into Dynamics of Mobile Genetic Elements in
24 Hyperthermophilic Environments from Five New Thermococcus Plasmids, *PLoS ONE*, **8**,
25 e49044.
- 26 Liu, L., *et al.* (2012) The human microbiome: a hot spot of microbial horizontal gene transfer,
27 *Genomics*, **100**, 265-270.
- 28 Pagani, I., *et al.* (2012) The Genomes OnLine Database (GOLD) v.4: status of genomic and
29 metagenomic projects and their associated metadata, *Nucleic Acids Res*, **40**, D571-579.
- 30 Rost, B. (2002) Enzyme function less conserved than anticipated, *J Mol Biol*, **318**, 595-608.
- 31 Saye, D.J., *et al.* (1987) Potential for transduction of plasmids in a natural freshwater
32 environment: effect of plasmid donor concentration and a natural microbial community on
33 transduction in *Pseudomonas aeruginosa*, *Applied and Environmental Microbiology*, **53**, 987-
34 995.
- 35 Sayers, E.W., *et al.* (2009) Database resources of the National Center for Biotechnology
36 Information, *Nucleic Acids Res*, **37**, D5-15.
- 37 Tatusov, R.L., *et al.* (2003) The COG database: an updated version includes eukaryotes,
38 *BMC Bioinformatics*, **4**, 41.

1 Zhu, C., *et al.* (2015) Functional Basis of Microorganism Classification, *PLoS Comput Biol*,
2 **11**, e1004472.

3

4