# Comparative systems biology across an evolutionary gradient within the *Shewanella* genus

Konstantinos T. Konstantinidis[1#*], Margrethe H. Serres[2#], Margaret F. Romine[3#], Jorge L. M. Rodrigues[4#], Jennifer Auchtung[5], Lee-Ann McCue[3], Mary S. Lipton[3], Anna Obraztsova[6], Carol S. Giometti[7], Kenneth H. Nealson[6], James K. Fredrickson[3], and James M. Tiedje[5*]

[#]Authors with equal contribution

Affiliations

[1]School of Civil and Environmental Engineering and School of Biology, Georgia Institute of Technology, Atlanta, Georgia 30332, USA.

[2]Josephine Bay Paul Center for Comparative Molecular Biology and Evolution, Marine Biological Laboratory, Woods Hole, Massachusetts 02543, USA.

[3]Biological Sciences Division, Pacific Northwest National Laboratory, Richland, Washington 99352, USA.

[4]Department of Biology, University of Texas, Arlington, Texas 76019, USA.
[5]Center for Microbial Ecology, Michigan State University, East Lansing, Michigan 48824, USA.
[6]Department of Earth Sciences, University of Southern California, Los Angeles, California 90089, USA.

[7]Biosciences Division, Argonne National Laboratory, Argonne, IL 60439.

Classification: BIOLOGICAL SCIENCES, Microbiology
Manuscript information: 19 pages; about 48,900 characters total (including characters displaced by figures)

*Authors for correspondence:
James M. Tiedje
Center for Microbial Ecology, 540 Plant and Soil Science Building, Michigan State University, East Lansing, MI 48824
Phone: (517) 353-7858. Fax: (517) 353-2917.
E-mail: tiedjej@msu.edu

Konstantinos T Konstantinidis
311 Ferst Drive, ES&T, Room 3224
Georgia Institute of Technology
Atlanta, GA 30332-0512
Phone: 404.385.3628. Fax:   404.894.8266
Email: kostas@ce.gatech.edu

1

**ABSTRACT**

To what extent genotypic differences translate to phenotypic variation remains a poorly understood issue of paramount importance for several cornerstone concepts of microbiology including the species definition. Here, we take advantage of the completed genomic sequences, expressed proteomic profiles, and physiological studies of ten closely related *Shewanella* strains and species to provide quantitative insights into this issue. Our analyses revealed that, despite extensive horizontal gene transfer within these genomes, the genotypic and phenotypic similarities among the organisms were generally predictable from their evolutionary relatedness. The power of the predictions depended on the degree of ecological specialization of the organisms evaluated. Using the gradient of evolutionary relatedness formed by these genomes, we were able to partly isolate the effect of ecology from that of evolutionary divergence and rank the different cellular functions in terms of their rates of evolution. Our ranking also revealed that whole-cell protein expression differences among these organisms when grown under identical conditions were relatively larger than differences at the genome level, suggesting that similarity in gene regulation and expression should constitute another important parameter for (new) species description. Collectively, our results provide important new information towards beginning a systems-level understanding of bacterial species and genera.

1     **\body**

2     **Introduction**

3        Predicting the phenotype of newly isolated organisms based upon the existing

4     knowledge of previously characterized organisms constitutes one of the most fundamental goals

5     of microbiology. Organisms isolated from diverse environments and habitats often have their

6     phenotypic and physiological properties inferred from their evolutionary relatedness, measured

7     by (mainly) the 16S rRNA gene sequence identity or other means (1, 2), to the type strains of

8     known species. Although this practice has been broadly applied in studies of microbial

9     communities, contributing greatly toward advancing microbiology knowledge, its use in this

10     manner is rooted in rather low-resolution experimental methods and procedures (1, 3). The

11     powerful genomic tools now available provide the opportunity for a much more detailed and

12     informative evaluation of the relationship between genetic and phenotypic similarity. Simple

13     questions that remain unanswered or only partially explored such as "to what degree do

14     microorganisms encode and express the same metabolic pathways when grown under identical

15     conditions" and "to what extent are the similarities in expressed pathways determined by the

16     genetic relatedness and/or the (distinct) ecological adaptations of the microorganisms?" can now

17     be answered accurately and quantitatively. Addressing such questions will provide long-needed

18     information to better understand and model the enormous microbial biodiversity that exists on

19     the planet.

20        To this end, we have analyzed and compared, both at the whole-genome and the whole-

21     proteome levels, ten isolates belonging to the genus *Shewanella*, an important genus in cycling

22     of organic and inorganic materials in the environment (4). These isolates originated from

23     diverse geographic locations and habitats, including fresh and marine water columns, sediments,

1  and subsurface environments (Fig. 1A and Table S1), and carry out a diverse range of metabolic

2  processes (4). Although precise ecological information, e.g., *in-situ* abundance and persistence

3  in time, about each isolate is typically not available, the procedure employed to isolate these

4  strains, i.e., enrichment cultures from a variety of environmental samples for the phenotype or

5  genotype of interest, is similar to common microbiology practice. Accordingly, our analyses

6  with the *Shewanella* strains should be relevant for the questions described above and for

7  broadening our understanding of the interrelationship between genotype, phenotype,

8  environment and evolution. Our results represent the first thorough and system-level assessment

9  of an environmental representative of *Proteobacteria*, an enormously diverse and important

10  group, that can be compared and contrasted to previous assessments of the heavily sampled

11  human pathogens or the ecologically specialized organisms such as the photosynthetic

12  *Prochlorococcus* (5). Such comparisons identified several trends that may apply to other

13  environmentally versatile bacteria besides *Shewanella*.

14

15  **A continuous genetic gradient within a genus.** Phylogenetic analysis of the 16S rRNA gene

16  sequences revealed that the ten *Shewanella* isolates formed a tight cluster, with the intra-cluster

17  sequence identity ranging from 92 to ~100% (Fig. 1B). Hence, these genomes belong justifiably

18  to the same genus according to the most frequently used standards of bacterial taxonomy (2, 6).

19  To gain further insights into the diversity of this group, the Average Nucleotide Identity (ANI)

20  of all pair-wise conserved genes between (any) two genomes, a more sensitive parameter for

21  measuring evolutionary relatedness among closely related genomes than the 16S rRNA gene

22  (7), was employed. The ANI analysis revealed that these genomes form a continuing gradient of

23  genetic relatedness, which was not readily apparent from the 16S rRNA gene analysis (Fig. 1C).

1    In particular, *S. putrefaciens* strains W3-18-1 and CN-32 as well as *Shewanella* sp. MR-4 and

2    MR-7 are the most closely related pairs, showing ANI values of ~96.5% and ~98.4%,

3    respectively. These values are well above the 95% ANI that corresponds to the 70% DNA-DNA

4    hybridization (DDH) standard frequently used for species demarcation, which is consistent with

5    the experimentally derived DDH values for these organisms (6). Hence, these pairs of genomes

6    sample the sub-species level. The MR-4 and MR-7 genomes show ~92%, ~85%, and ~79%

7    ANI to *Shewanella* sp. ANA-3, *S. putrefaciens* CN-32, and *S. oneidensis* MR-1 genomes,

8    respectively. Thus, these genome pairs represent varied levels of genetic relatedness within the

9    *Shewanella* genus. Finally, all the previously mentioned genomes show ~69.7-72% ANI to *S.*

10    *frigidimarina* NCIMB400, *S. denitrificans* OS217, *S. loihica* PV-4, and *S. amazonensis* SB2B

11    strains, which represent the four most divergent species sampled within the genus. This gradient

12    provided the opportunity to precisely estimate the number of changes in the genes, pathways

13    and subsystems of the cell over time and as a result of environmental adaptations and selection

14    pressures.

15

16    **Gene-content variation as a function of evolutionary time and ecology.** The ten *Shewanella*

17    isolates have similar genome sizes, varying from 4.3 to 5.3 Mbp (Table S1). Comparative

18    analysis revealed extensive gene content diversity among the genomes. From the 9,782

19    predicted non-redundant (orthologous genes removed) protein-coding sequences (CDS)

20    annotated in the ten genomic sequences (the pangenome) only ~2,128 (22%, constituting ~54%

21    of the total genes in the genome, on average) were present in all genomes (core CDS set); about

22    2,965 (30%) were found in at least two genomes (variable CDSs), while the remaining CDS

23    (4,689 or 48%) were strain specific (Fig. 2B and Table S2). Nonetheless, the majority of the

1    variable CDSs were found to be specific to clades, i.e., the MR or *S. putrefaciens* clades (Fig.

2    1B), while a smaller fraction had a more sporadic distribution among the strains (see the

3    similarity between the gene content tree and the phylogeny of the genomes in Fig. 3).

4    Accordingly, the overall extent of CDS-content similarity showed a very strong linear decrease

5    with increasing evolutionary distance between the genomes compared ($R^2 > 0.9$, see Fig 1C),

6    which is consistent with results reported previously based on other bacterial groups (7). The

7    strong linear trend suggests that, despite the extensive gene diversity and apparent genome

8    fluidity, the genotypic similarity of bacteria may be generally predictable from their

9    evolutionary relatedness.

10          Although a tight relationship between shared CDS-content and evolutionary relatedness

11   was observed, several significant departures (outliers) from this main trend were also noted and

12   were most likely attributable to ecological adaptations. For instance, the two most closely

13   related genomes based on ANI, CN32 and W3-18-1 (98.4% ANI), showed substantially more

14   CDS-content differences compared to what was expected based on their small evolutionary

15   divergence (see regression trendline in Fig. 1C) or compared to the more distantly related

16   (96.4% ANI) pair of MR-7 and MR-4 (~530 vs. ~430 CDSs, respectively, not counting CDS on

17   mobile elements; see Table S2). CN-32 and W3-18-1 were isolated from more diverse

18   environments (deep-subsurface sandstone vs. marine sediment, respectively) compared to MR-4

19   and MR-7 (5m vs. 60m depth in the Black Sea, respectively). Hence, it is likely that genetic

20   adaptations specific to these environments account for the larger gene content differences

21   observed in the former strains relative to the latter ones. In agreement with the latter

22   interpretation, CN-32-specific genes included several genes that might be important for survival

1    in the subsurface environment such as an arsenate reductase, copper resistance system, heavy

2    metal efflux pump, and a polysaccharide biosynthesis cluster.

3           Similarly, *S. denitrificans* strain OS217 is as divergent as three other isolates (strains

4    PV-4, NCIMB400, and SB2B) are from the remaining six *Shewanella* isolates in our collection

5    (e.g., Fig. 3D). Yet, the OS217 genome contained substantially more strain-specific genes and

6    showed the greatest loss of "core-like" CDSs (i.e., CDSs present in all other *Shewanella*

7    genomes) compared to the genomes of PV-4, NCIMB400, or SB2B (Table S2). For instance,

8    the core set increased by 265 genes when OS217 was removed from the analysis compared to

9    fewer than 60 genes when PV-4, NCIMB400 or and SB2B were individually removed. Our

10    genomic, physiological (e.g., Table S5), and proteomic data collectively suggests that strain

11    OS217 has undertaken a unique evolutionary path, possibly driven by the loss of the three

12    menaquinone biosynthetic gene clusters (*menDHCE*, *menF*, *menB*) common to the other

13    *Shewanella* strains and resulting in inability to exploit strictly anaerobic habitats. These results

14    are also consistent with previous findings suggesting that strain OS217 is a specialized

15    denitrifier (4) and with the longstanding observation that respiratory denitrification is not found

16    in organisms that are strong fermentors (8). These findings may indicate that more extensive

17    genetic changes are involved for an organism to diverge to the opposite physiology. Lastly, the

18    (outlier) pairs of genomes with a higher percentage of shared genes than the average, i.e., CN32

19    or W3-18-1 vs. MR4 or MR7 (Fig. 1C), are attributable to the substantially smaller size of these

20    genomes (i.e., 4.6-4.7 Mbs) relative to that of the rest of the genomes (i.e., ~5.2 Mbp, see Table

21    S1) rather than to more similar ecological adaptations (the number of shared orthologs and

22    mobile gene content in these pairs is comparable to that of other pairs).

23

1   **Processes contributing to gene-content variation.** To provide further quantitative insights into

2   the processes contributing to gene-content variation, the genes that differed in pair-wise whole-

3   genome comparisons were assigned to five major functional categories and the percentage of

4   genes in each category was evaluated against the genetic relatedness of the two genomes

5   compared. The five categories were: i) pseudogenes, denoting genes predicted to encode

6   insertions, deletions, or sequence alterations that would result in premature termination of the

7   encoded protein, ii) IS/Tn, denoting insertion sequences or transposons, iii) mobile islands,

8   denoting runs of neighboring genes (genomic islands) that included integrase genes, iv) other,

9   denoting all other unique genes, including genomic islands that do not contain clear evidence of

10   being mobile, and v) hypothetical or conserved hypothetical, denoting the fraction of the genes

11   in category (iv) that had no detectable homolog in any of the fully sequenced genomes except in

12   other *Shewanella* genomes (Table S3). Our results revealed that mobile islands and insertion

13   elements dominated the gene-content differences among genomes of the same species but their

14   contribution gradually decreased in comparisons among genomes of increasing evolutionary

15   divergence at the expense of genes in the "other" category (Fig. 2A). These findings are

16   consistent with rapid turnover of mobile islands over short evolutionary scales. Further, the

17   majority (>75%) of the genes in the "other" category were typically found in clusters of ~5 to

18   ~40 genes, reflecting presumably their "mobile island" origin. These findings are consistent

19   with preferential deletion of the mobility/transposition genes (presumably due to negative

20   selection) in the course of evolution and retention of only the potentially ecologically important

21   genes of mobile islands. Therefore, the *Shewanella* organisms evaluated here appear to acquire

22   most of their new functions as follows: acquisition of mobile islands followed by selection for

1 the islands carrying ecologically important genes and finally loss of the mobile and ecologically

2 unimportant genes.

3

4 **The *Shewanella* pangenome and conserved gene core.** Comparative analysis of the ten

5 *Shewanella* genomic sequences revealed that sampling of the genus pangenome remained

6 unsaturated (Fig. 2B, blue bars); this result was attributable to the large number (468, on

7 average) of strain-specific genes. Only 10% to 25% of the latter genes, depending on the

8 genome evaluated, found a homolog in a genome outside the *Shewanella* genus when queried

9 against all bacterial genomes available at the end of 2008, indicating the great potential for

10 discovering novel genes with more *Shewanella* strains sequenced. The number of new genes per

11 genome is an order of magnitude higher than those calculated for highly specialized human

12 pathogens (9) but significantly lower than that of the opportunistic pathogen *Escherichia coli*

13 *(7)*. It must be pointed out, however, that these pan-genome calculations are not directly

14 comparable and should be interpreted with caution. For instance, the average ANI value among

15 all pairs of *Shewanella* genomes is ~76%, which is significantly lower than that within the *E.*

16 *coli* group (~96%), and there appears to be a strong positive correlation between the amount of

17 novel genes carried in a genome and the (higher) degree of evolutionary divergence of the

18 genome, regardless of the effect of ecology or environmental adaptation (Fig. 1C and in (7)). On

19 the other hand, the prophage content of the *E. coli* genomes is substantially higher than that of

20 the *Shewanella* ones (10-20% vs. 0-5%, respectively), and this accounts for much of the

21 difference observed. When the groups were adjusted for comparable intra-group diversity, by

22 including selected *Salmonella* (~82% ANI to *E. coli*) and *Yersinia* (~72% ANI to *E. coli*)

23 genomes together with *E. coli* ones and with prophage genomes removed from the analysis, the

1  gene diversity observed within the enterics was comparable to that of the *Shewanella* (Fig. 2B).

2  Therefore, the evaluation of these two important groups suggests that sequencing of any new

3  organism, as long as the organism belongs to a versatile genus and has a different ecological

4  history relative to the previously sequenced members of the genus, should be expected to

5  expand substantially the pangenome of the genus.

6  Both the *Shewanella* and the enterics core gene sets were highly enriched in

7  translational, transcriptional, DNA replication, and central metabolism genes and overlapped

8  extensively (~50% of the genes were shared between the two cores). *Shewanella*-specific core

9  functions were associated mainly with metabolic pathways, as well as chemotaxis and sensory-

10  transduction processes. Using the BioCyc pathway schema (10), 104 pathways were identified

11  as being common to all *Shewanella* genomes, including pathways for energy metabolism,

12  synthesis of building blocks (amino acids, cofactors, fatty acids, and nucleotides), and for

13  degradation or inter-conversion of metabolites and all but two amino acids and metabolites (Fig.

14  S2, and Table S4). A common trait of the *Shewanella* strains appears to be the use of the

15  pentose phosphate and Entner-Doudoroff pathways for hexose degradation. This is based on the

16  lack of the enzyme 6-phosphofructokinase (Pfk; the most important regulatory enzyme of the

17  canonical glycolysis pathway), initially observed in previous gene expression studies of MR-1

18  cultures (11). Members of the *Shewanella* genus also have fewer phosphotransferase system

19  (PTS) transporters than usually encountered in proteobacterial genomes. Whether there is a

20  connection between the reduced PTS and lack of Pfk is not clear, but it is possible that the lower

21  level of phosphoenolpyruvate (PEP) synthesized as a result of not using the glycolytic pathway

22  may render the PEP-dependent PTS system inefficient.

1      When the core was defined as the genes present in all but one of the 10 genomes, the

2    dataset increased by 411 protein coding genes (265 when OS217 was excluded from the

3    analysis), corresponding to, on average, 12-14% of the *Shewanella* genome (Table S2). These

4    findings suggest that gene loss, including loss of genes that are apparently indispensable for the

5    majority of the strains of a species, might be a successful strategy for fast evolution and

6    environmental adaptation. A representative example of strain-specific adaptations related to a

7    group core function, which involved considerable gene deletion and/or gene acquisition, is

8    given below. All *Shewanella* strains except for *S. denitrificans* OS217, which shows limited

9    anaerobic growth capabilities presumably due to gene loss during the process of ecological

10   specialization (discussed above), were able to reduce several metals and metalloids (Table S5),

11   a well-known characteristic of the genus (12). The main metal reductase locus, encoded by

12   *mtrCAB* genes, is virtually identical for the nine strains but the adjacent loci vary, reflecting

13   evolutionary history and possibly metal respiratory specialization (4). These dissimilarities

14   explained some, but not all, of the variation in metal respiration among strains observed during

15   our growth experiments. For example, although their *mtr* locus and flanking genes are identical,

16   strain CN-32 was able to grow on lactate (20 mM) when six different metals or metalloids were

17   used as electron acceptors, whereas strain W3-18-1 only grew with Fe, Mn, and Se, under the

18   conditions tested. These results may reflect differences in the upstream pathways to metal

19   reduction between the two strains and underscore the need for more research to understand

20   better the details of the metal respiration cascade.

21

22   **Gene presence vs. expression as a function of time and ecology.** Transcriptome comparisons

23   have shown that gene expression rather than gene content differences, occurring either at

1    different times and/or tissues, are mainly responsible for the differential development of

2    eukaryotic organisms, e.g. human and chimpanzees (13), and the adaptive evolution of natural

3    populations (14). It follows that, in addition to the number of shared genes, gene expression

4    constitutes an important factor determining phenotypic similarity (or dissimilarity). While the

5    latter applies presumably to bacteria as well, systematic assessments of the role of gene

6    expression on the phenotypic differences observed among closely related organisms are lacking.

7    To begin exploring this issue, the ten strains were grown under identical batch-culture

8    conditions to obtain their whole-cell proteome profiles and contrast the profiles against the

9    evolutionary relatedness among the strains. Overall, the degree of similarity in proteome

10   profiles was congruent with the evolutionary relatedness among the strains, i.e., the fraction of

11   orthologous proteins detected to be expressed in the cultures was higher in closely related

12   strains than in more divergent strains. However, the differences in expressed proteins among the

13   strains were consistently larger than their differences at the gene-content level when gene

14   expression and gene content were assessed for the same 4,300 (reference) genes found in the

15   MR-1 genome (compare branch lengths in Fig. 3B and 3C), which minimized the effect of

16   gene- or strain- specific variations in the measurements. More surprisingly, the same pattern

17   was observed even when gene expression was assessed for the core genes only (Fig. 3A, Fig.

18   S2), which circumvented the dependency of the proteome profiles on the underlying gene-

19   content differences in the previous comparisons. These results were attributable to a high

20   number of proteins expressed by one or a few, but not all, of the strains possessing the

21   corresponding gene, with proportions that varied from 1.9 to 2.6 times more than those proteins

22   expressed by all strains possessing the corresponding gene (Table S6). For instance, although

23   twenty percent (556 genes) of the core proteins were expressed by all strains, a substantially

1   larger fraction of core proteins (993 or 36%) were expressed by one or more (but not all) strains.

2   While some of these differences may be due to higher noise in the proteomics data relative to

3   the genomics data, we believe that many of these differences are biologically relevant due to the

4   high reproducibility (>80%) of proteomics measurements on batch cultures like the ones used in

5   the present study (15), our high stringency in processing and analyzing the proteomics data (see

6   methods), and the fact that very similar results were found when a subset of five specific

7   regions of traditional 2-dimension protein gels were overlaid and compared for absence or

8   presence of protein spots (Fig. S3). Finally, proteins characteristic of the stationary growth

9   phase, such as the RpoS sigma factor (16), were not detected in the expressed proteomes,

10  suggesting that all of our cultures were sampled at their exponential growth phase.

11         Our findings revealed that although strains CN-32 and W3-18-1 are significantly more

12  closely related than are strains MR-4 and MR-7 [e.g., a 2% higher ANI value translates to

13  substantially higher gene-content and evolutionary relatedness, as we and others have shown

14  (7)], they showed comparable differences in expressed proteins compared to the latter strains for

15  the same genes analyzed (Fig. 3). These findings could therefore be attributable to a higher

16  degree of environmental/ecological adaptations (which may have altered metabolic and

17  regulatory networks) in the CN-32/W3-18-1 pair relative to the MR-4/MR-7 one. Similarly, *S.*

18  *denitrificans* OS217, which appeared to be the most ecologically specialized organism of the

19  set, also showed the most unique proteomic profile (Fig. 3). The larger gene expression

20  differences observed for OS217 and CN-32/W3-18-1 than anticipated based on their

21  evolutionary divergence alone echoes the results described above based on the gene-content

22  analysis. Further, the largest fraction (44%) of the proteins detected in the protein profiles was

23  strain-specific and included many non-hypothetical proteins such as outer membrane proteins,

TonB-dependent receptors, proteases, restriction-modification enzymes, glycosylases, and polysaccharide biosynthesis enzymes. Most of these proteins can be linked to metabolic fitness or interaction with the environment, and hence could possibly underlie important physiological and/or regulatory differences among the strains. The extensive variability in core proteins and the high number of strain-specific proteins expressed under identical growth conditions indicates a multifaceted and highly dynamic control of whole genome expression. Collectively, our proteomics analyses suggest that changing this control appears to represent a particularly important mechanism, in addition to gene acquisition or loss, for fast adaptation in changing and diverse environments. Consistent with these conclusions, the first mutations observed in experimentally evolved *E. coli* strains for 20,000 generations under laboratory conditions involved regulatory genes and networks (17).

**Compartmentalized microbial evolution.** In order to characterize which cellular functions evolve faster in the *Shewanellae*, the percent conservation of selected functional gene categories (see methods for details) was evaluated against the evolutionary relatedness among the strains compared (measured by % ANI). As evolutionary distance increased, the % conservation of all categories decreased, but the extent of decline (i.e., the slope) differed, presumably reflecting the varied selection pressures on the corresponding genes. The analysis revealed the following order: pathways were substantially more conserved than individual orthologs, orthologs more conserved than transcriptional regulators, sensing and respiration genes, and expressed proteins (Fig. 4). The most rapidly changing individual functions, both in terms of gene presence/absence and sequence conservation, were TonB-dependent outer membrane receptors followed by methyl-accepting chemotaxis proteins, transcription regulators and cytochromes.

1   These results are consistent with our previous findings and suggest that genomic and regulatory

2   changes in sensing mechanisms represent the first line of adaptive response to different redox

3   conditions. Experimentally determined anaerobic growth characteristics such as biomass

4   produced and electron acceptors used (Table S5) were also very different among the *Shewanella*

5   strains and ranked among the fastest changing functional entities (Fig. 4). A growth phenotype

6   encompasses the sensing of a substrate, expression of relevant regulators, transporters and

7   enzymes, in addition to physiological parameters related to the change in growth conditions.

8   These potential sources for additional variation among the strains may explain why the growth

9   phenotype is significantly less conserved compared to pathways, orthologs, and protein

10  expression patterns.

11

12  **Summary and perspectives for the future**

13       Microbiologists have been primarily focused on comparisons among either very closely

14  related strains of the same species or distantly related species in order to advance understanding

15  of the microbial life on Earth. The ten *Shewanella* genomes studied here were selected to

16  represent a range of evolutionary distances, providing for a more unconstrained view of

17  microbial diversity and evolution. Comparisons among these genomes revealed that the

18  *Shewanella* genus is genomically and more so proteomically diverse. Although a high degree of

19  variation in protein expression profiles was anticipated among distantly related species, the

20  variation observed among strains of the same species was comparatively much larger than

21  expected, given also the single growth condition used (Fig. 3 & 4). It also appears that, in some

22  cases, the variation in expressed proteomes correlated positively with the extent of

23  environmental adaptation (specialization). These findings have important implications for the

1  correspondence between genotype and phenotype and hence, for the bacterial species concept.

2  The evolutionary and functional gradients reported here also suggested that specialization might

3  occur over a very short time span, much shorter compared to what corresponds to the current

4  species standards. Specialization appeared to take place primarily through changes at the

5  regulatory level and through the high plasticity and fluidity of the *Shewanella* chromosomes

6  (e.g., Fig. 4).

7      The power of "omics" compared to traditional approaches to unravel organism's

8  environmental/ecological adaptations and make robust predictions about the similarity (or

9  difference) in phenotypic traits among organisms was also highlighted by our analyses. The

10  literature as well as our experimentally derived physiological and growth data could not easily

11  distinguish between most of the strains used in this study or (even) define general properties for

12  the major clades represented by these strains. This was also reflected in the very low correlation

13  obtained between anaerobic growth characteristics (Table S5) and the evolutionary relatedness

14  of the strains compared. In contrast, genomic and proteomic data correlated well with the

15  phylogeny of the strains and identified congruently strain-specific adaptations that might be

16  linked to speciation for several of the strains studied. These results further corroborate the

17  notion that it is time to start replacing the traditional approaches for defining diagnostic

18  phenotypes for (new) species or clades with omics-based procedures.

19      Distinguishing the effect of ecological adaptation from that of evolutionary divergence

20  alone represents the most limiting factor in increasing the power of our predictions on

21  phenotype based on the genotype. Towards this direction, studying the extent of variation

22  among members of the same natural population, i.e., among organisms with very similar

23  environmental adaptations, and contrasting it to the levels of variation detected in this study

1   with diverse organisms will allow for fruitful conclusions. The trendlines obtained in this study

2   (Fig. 1C & Fig. 4) also provide a reference for comparing organisms of narrower (or broader)

3   metabolic versatility than the *Shewanellae*. Further, although the growth conditions used in this

4   study were very limited, they remain artificial compared to the environmentally relevant

5   conditions and hence, may represent different stresses for each isolate evaluated. Replicate

6   experiments and experiments performed with continuous cultures (chemostats) are currently

7   underway in order to provide further quantitative insights into the role of variation in gene

8   expression. Finally, a major limitation remains in that, despite the dedicated efforts of numerous

9   laboratories, many of the genes in the genome have not been experimentally characterized and

10  their physiological role is unknown. Continuing the efforts to establish function to as many

11  genes in the genome as possible is critical for a thorough understanding of a bacterium that

12  could serve as a model for versatile environmental bacteria.

13      Regardless of these limitations, the results presented here constitute important

14  information towards better modeling the correspondence between genotype and phenotype and

15  provide directions and testable hypotheses that will bring us one step closer to systems-level

16  understanding of microbial species and populations.

17

18  **Material and Methods**

19  The organisms used in this study, their genomic features, gene-content, and accession numbers

20  of the versions of the genomic sequences used in the study are provided in Table S1. Orthologs

21  were identified for the ten *Shewanella* genomes by a combination of three methods: i) protein-

22  protein pair-wise reciprocal BLAST (blastp) (18), ii) reciprocal protein-genomic sequence best

23  match (tblastn), and iii) Darwin pair-wise best hit (19). Genes found in plasmids or mobile

1    elements were excluded from ortholog and proteome comparisons among the strains. The

2    degree of conservation of cellular functions or traits between two strains (e.g., Fig. 4) was

3    determined as follows. I) For orthologs, transcriptional regulators, TonB receptors, MCPs, and

4    cytochromes: all genes in the genome assignable to each of these categories were determined

5    based on the gene annotation and the number of orthologous genes shared between two strains

6    for each category (according to Table S2) was divided by the total number of genes assignable

7    to the category for each strain. The two values were averaged to provide the values used in

8    figure 4. II) A total of 163 unique pathways were identified in the ten *Shewanella* genomes

9    according to the BioCyc pathway schema (http:/biocyc.org). The number of shared pathways

10    between the strains, as a fraction of the total pathways carried by a strain, was determined based

11    on the presence/absence of the corresponding pathway genes. III) For proteomes and anaerobic

12    growth, the number of orthologous proteins expressed (Table S6) and metal/metalloids respired

13    (Table S5) by both strains in a pair was divided by the total number of (non-redundant) proteins

14    expressed and metal/metalloids respired by either strain, respectively. The use of "total traits

15    counted for both strains" as the denominator (as opposed to "counts for one strain") provided

16    also for more direct comparisons to the sequence-based traits (I and II above) because otherwise

17    the latter traits would have been penalized relatively higher due to the high number of

18    "auxiliary" genes, which remained un-expressed under the simple growth conditions tested. For

19    proteomics analysis, cultures were grown aerobically in Tryptic Soy Broth to final Optical

20    Density, OD=0.5. Cells were lysed, proteins extracted and digested with trypsin, and the

21    resulting peptides analyzed by mass spectrometry as previously described (20), with the only

22    exception that filtering of the data was performed as described in (21). Two-dimensional

1 proteomic gels were carried out as described previously (15). A detailed description of materials

2 and methods is included in the supplementary material.

3

4

5

6 **Acknowledgements**

**REFERENCES**

1.      Stackebrandt E*, et al.* (2002) Report of the ad hoc committee for the re-evaluation of the species definition in bacteriology. *Int J Syst Evol Microbiol* 52(3):1043-1047.
2.      Brenner D, Staley J, & Krieg N (2001) *Bergey's manual of systematic bacteriology* (Springer-Verlag, New York) 2nd Ed pp 27-31
3.      Vandamme P*, et al.* (1996) Polyphasic taxonomy, a consensus approach to bacterial systematics. *Microbiol Rev* 60(2):407-438.
4.      Fredrickson JK*, et al.* (2008) Towards environmental systems biology of Shewanella. *Nat Rev Microbiol* 6(8):592-603.
5.      Kettler GC*, et al.* (2007) Patterns and implications of gene gain and loss in the evolution of Prochlorococcus. *PLoS Genet* 3(12):e231.
6.      Goris J*, et al.* (2007) DNA-DNA hybridization values and their relationship to whole-genome sequence similarities. *Int J Syst Evol Microbiol* 57(Pt 1):81-91.
7.      Konstantinidis KT, Ramette A, & Tiedje JM (2006) The bacterial species definition in the genomic era. *Philos Trans R Soc Lond B Biol Sci* 361(1475):1929-1940.
8.      Tiedje JM (1988) *Ecology of denitrification and dissimilatory nitrate reduction to ammonium* (John Wiley and Sons, New York) pp 179-244.
9.      Medini D, Donati C, Tettelin H, Masignani V, & Rappuoli R (2005) The microbial pan-genome. *Curr Opin Genet Dev* 15(6):589-594.
10.     Caspi R*, et al.* (2008) The MetaCyc Database of metabolic pathways and enzymes and the BioCyc collection of Pathway/Genome Databases. *Nucleic Acids Res* 36(Database issue):D623-631.
11.     Driscoll ME*, et al.* (2007) Identification of diverse carbon utilization pathways in Shewanella oneidensis MR-1 via expression profiling. *Genome Inform* 18:287-298.
12.     Hau HH & Gralnick JA (2007) Ecology and biotechnology of the genus Shewanella. *Annu Rev Microbiol* 61:237-258.
13.     Enard W*, et al.* (2002) Intra- and interspecific variation in primate gene expression patterns. *Science* 296(5566):340-343.
14.     Oleksiak MF, Churchill GA, & Crawford DL (2002) Variation in gene expression within and among natural populations. *Nat Genet* 32(2):261-266.
15.     Elias DA*, et al.* (2008) The influence of cultivation methods on Shewanella oneidensis physiology and proteome expression. *Arch Microbiol* 189(4):313-324.
16.     Lange R & Hengge-Aronis R (1991) Identification of a central regulator of stationary phase gene expression in Escherichia coli. . *Mol. Microbiol.* (5):49-59.
17.     Philippe N, Crozat E, Lenski RE, & Schneider D (2007) Evolution of global regulatory networks during a long-term experiment with Escherichia coli. *Bioessays* 29(9):846-860.
18.     Altschul SF*, et al.* (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 25(17):3389-3402.
19.     Gonnet GH, Hallett MT, Korostensky C, & Bernardin L (2000) Darwin v. 2.0: an interpreted computer language for the biosciences. *Bioinformatics* 16(2):101-103.

20. Fang R*, et al.* (2006) Differential label-free quantitative proteomic analysis of Shewanella oneidensis cultured under aerobic and suboxic conditions by accurate mass and time tag approach. *Mol Cell Proteomics* 5(4):714-725.

21. Washburn MP, Wolters D, & Yates JR, 3rd (2001) Large-scale analysis of the yeast proteome by multidimensional protein identification technology. *Nat Biotechnol* 19(3):242-247.

22. Bruen TC, Philippe H, & Bryant D (2006) A simple and robust statistical test for detecting the presence of recombination. *Genetics* 172(4):2665-2681.
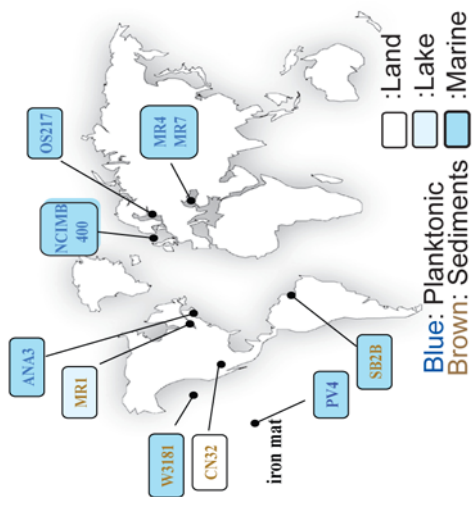
**FIGURE LEGENDS**

**Figure 1. The ten *Shewanella* genomes used in this study and their evolutionary gradient.** The geographic origin (**A**) and the 16S rRNA-based phylogenetic tree (**B**) of the ten genomes (in bold) are shown. The scale represents the number of substitutions per position and the numbers above and below the nodes represent the bootstrap support from 1,000 re-samplings using parsimony and maximum likelihood methods, respectively. Bootstrap values below 50 were omitted. A continuous genetic gradient was formed (**C**) when the fraction of the total genes in the genome shared between two genomes (y-axis) was plotted against the ANI of the shared genes between the two genomes (45 comparisons in total are shown). Dashed blue lines represent the 90% prediction intervals of the regression line; thus, open squares identify the outlier pairs of genomes observed (discussed in the text).

**Figure 2. The *Shewanella* pangenome.** *A: Contribution of different categories of genes to the pangenome as a function of ANI*. The genes that differed in all pair-wise whole-genome comparisons among the ten *Shewanella* genomes (45 comparisons in total) were assigned to five major functional categories (graph legend). The number of genes in each category, expressed as a fraction of the total genes that differed between the two genomes (y-axis), is plotted against the genomic ANI value of the two genomes compared. Individual data-points representing each comparison have been removed for clarity; only trendlines representing the mean and bars representing one standard deviation from the mean are shown instead. *B: Comparisons to the enterics pangenome*. The number of genes that remained conserved (y-axis) with the inclusion of more genomes in the analysis is plotted against the number of genomes (x-axis) used (light colors). The total number of non-redundant unique genes in all genomes used is also shown (dark colors). Bars represent one standard deviation based on 10 random combinations in adding the genomes to the analysis.
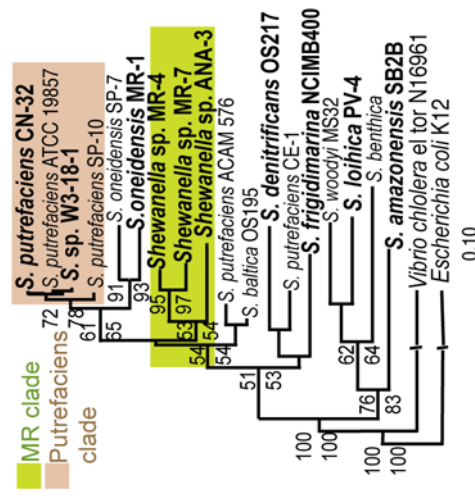
**Figure 3. Genome vs. proteome comparisons among nine *Shewanella* strains.** The protein profiles of nine *Shewanella* strains were compared based on the 2,128 core genes (**Panel A**) and the 4,300 genes found in the genome of strain MR-1 (**Panel B**) for gene expression, and the nine strains were subsequently clustered based on their overall similarity in the expression patterns of these two gene sets as follows: For each gene set, a full (all genes by all genomes) 0/1 matrix was built, with 1 denoting expression (defined as the detection of at least 2 unique peptides per protein) and 0 denoting no expression of the corresponding protein; the derived matrices were clustered as described in the supplementary material and the resulting cladograms are shown. Similarly, the nine strains were also clustered based on the presence/absence of the 4,300 MR-1 gene orthologs in their genome (sequence comparisons, **Panel C**). A maximum likelihood phylogenetic tree of the concatenated alignment of 1,507 single-copy core genes that had no detectable signal for recombination by Phi Test analysis (22) is also shown (**Panel D**). Scale bars represent percent similarity in the derived matrices for panels A, B, and C; and number of substitution per site for panel D.

**Figure 4. Modeling bacterial genotypic and phenotypic conservation across an evolutionary gradient.** The presence of orthologous proteins, TonB outer membrane receptors, cytochromes, methyl-accepting chemotaxis proteins (MCPs), transcriptional
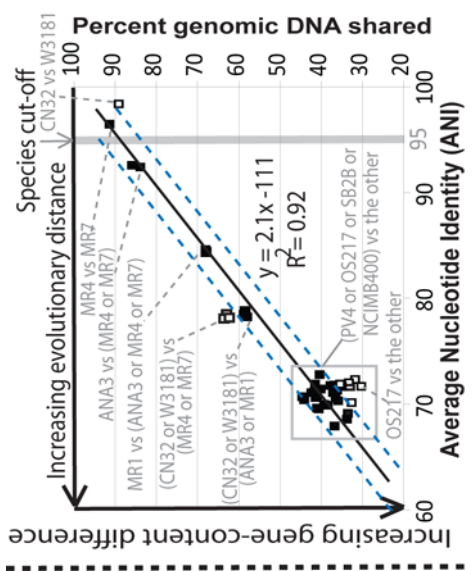
regulators, metabolic pathways, protein expression patterns, and reduction of metal or metalloids (anaerobic growth) was determined for the ten *Shewanella* strains (see methods). Each of the traits was compared among the *Shewanella* strains in a pair-wise manner (45 comparisons in total). The fraction of shared traits was determined for each pair of strains and plotted against the average nucleotide identity (ANI) of the respective strain pair. The inserted graph depicts the relationships between conservation of the traits and evolutionary distance using linear regression trendlines adjusted to intersect with the x and y-axis at 100%. The r-squared values of the regressions are also shown (figure legend).
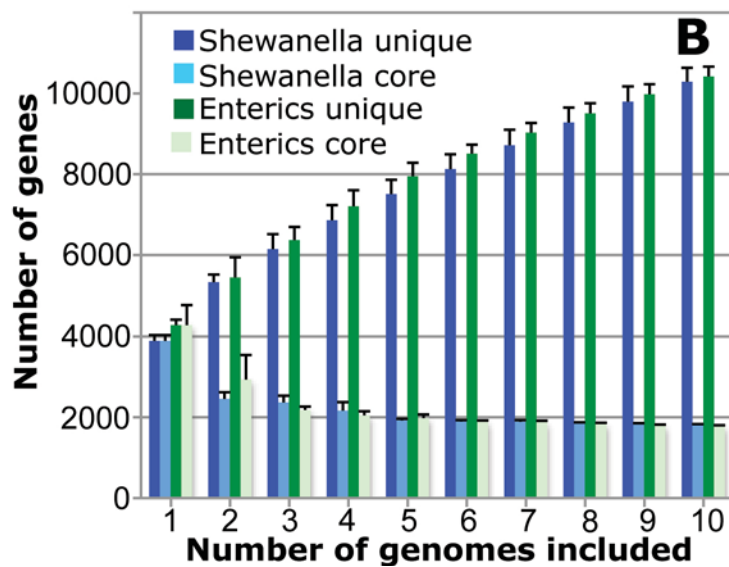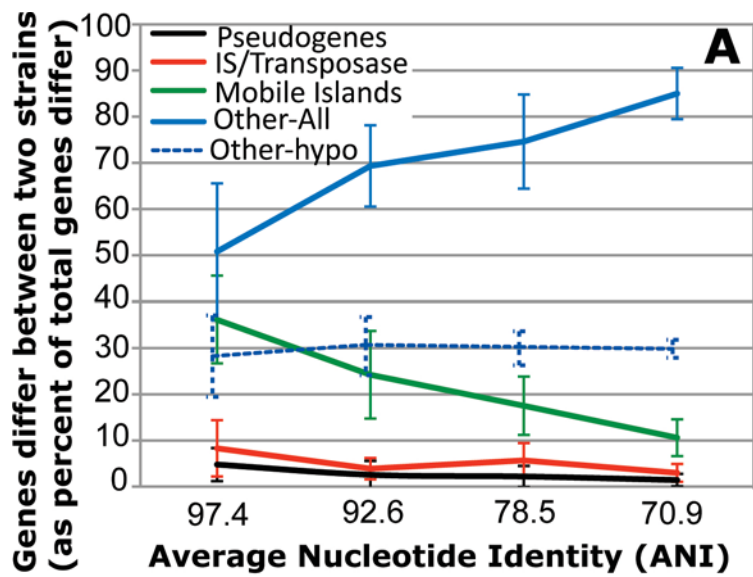
## C. An evolutionary gradient

Percent genomic DNA shared

$$y = 2.1x - 111$$
$$R^2 = 0.92$$

Species cut-off

Increasing evolutionary distance

Increasing gene-content difference

Average Nucleotide Identity (ANI)

CN32 vs W3181

MR4 vs MR7

ANA3 vs (MR4 or MR7)

MR1 vs (ANA3 or MR4 or MR7)

(CN32 or W3181) vs (MR4 or MR7)

(CN32 or W3181) vs (ANA3 or MR1)

(PV4 or OS217 or SB2B or NCIMB400) vs the other

OS217 vs the other

## B. The 16S rDNA tree

S. putrefaciens CN-32
S. putrefaciens ATCC 19857
S. sp. W3-18-1
S. putrefaciens SP-10
S. oneidensis SP-7
S. oneidensis MR-1
Shewanella sp. MR-4
Shewanella sp. MR-7
Shewanella sp. ANA-3
S. putrefaciens ACAM 576
S. baltica OS195
S. denitrificans OS217
S. putrefaciens CE-1
S. frigidimarina NCIMB400
S. woodyi MS32
S. loihica PV-4
S. benthica
S. amazonensis SB2B
Vibrio chlolera el tor N16961
Escherichia coli K12

MR clade
Putrefaciens clade

0.10

72
78
61
55
91
93
95
97
53
54
54
54
51
53
62
64
76
83
100
100
100
100

## A. The Geographic origin

OS217
MR4 MR7
NCIMB 400
ANA3
MR1
W3181
CN32
iron mat
PV4
SB2B

:Land
:Lake
:Marine

Blue: Planktonic
Brown: Sediments

Planktonic
Sediments

# PROTEOME CLUSTERING
## A: Core genes only

S. denitrificans OS217
S. frigidimarina NCIMB400
S. loihica PV-4
S. amazonensis SB2B
S. sp. W3-18-1
S. putrefaciens CN-32
Shewanella sp. ANA-3
Shewanella sp. MR-4
Shewanella sp. MR-7

0.1

## B: All MR-1 genes

S. denitrificans OS217
S. frigidimarina NCIMB400
S. loihica PV-4
S. amazonensis SB2B
S. sp. W3-18-1
S. putrefaciens CN-32
Shewanella sp. ANA-3
Shewanella sp. MR-4
Shewanella sp. MR-7

0.1

# C: GENE-CONTENT CLUSTERING

S. denitrificans OS217
S. frigidimarina NCIMB400
S. loihica PV-4
S. amazonensis SB2B
S. sp. W3-18-1
S. putrefaciens CN-32
Shewanella sp. ANA-3
Shewanella sp. MR-4
Shewanella sp. MR-7

0.1

# D: GENOME PHYLOGENY

S. loihica PV-4
S. amazonensis SB2B
S. denitrificans OS217
S. frigidimarina NCIMB400
S. putrefaciens CN-32
S. sp. W3-18-1
Shewanella sp. MR-4
Shewanella sp. MR-7
Shewanella sp. ANA-3

0.1