

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37

## Deep Annotation of Protein Function across Diverse Bacteria from Mutant Phenotypes

Morgan N. Price<sup>1</sup>, Kelly M. Wetmore<sup>1</sup>, R. Jordan Waters<sup>2</sup>, Mark Callaghan<sup>1</sup>, Jayashree Ray<sup>1</sup>, Jennifer V. Kuehl<sup>1</sup>, Ryan A. Melnyk<sup>1</sup>, Jacob S. Lamson<sup>1</sup>, Yumi Suh<sup>1</sup>, Zuelma Esquivel<sup>1</sup>, Harini Sadeeshkumar<sup>1</sup>, Romy Chakraborty<sup>3</sup>, Benjamin E. Rubin<sup>4</sup>, James Bristow<sup>2</sup>, Matthew J. Blow<sup>2,\*</sup>, Adam P. Arkin<sup>1,5,\*</sup>, Adam M. Deutschbauer<sup>1,\*</sup>

<sup>1</sup>Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory

<sup>2</sup>Joint Genome Institute, Lawrence Berkeley National Laboratory

<sup>3</sup>Climate and Ecosystem Sciences Division, Lawrence Berkeley National Laboratory

<sup>4</sup>Division of Biological Sciences, University of California, San Diego

<sup>5</sup>Department of Bioengineering, University of California, Berkeley

\*To whom correspondence should be addressed:

MJB ([MJBlow@lbl.gov](mailto:MJBlow@lbl.gov))

APA ([APArkin@lbl.gov](mailto:APArkin@lbl.gov))

AMD ([AMDeutschbauer@lbl.gov](mailto:AMDeutschbauer@lbl.gov))

Website for interactive analysis of mutant fitness data:

<http://fit.genomics.lbl.gov/>

Website with supplementary information and bulk data downloads:

<http://genomics.lbl.gov/supplemental/bigfit/>

38

39

## 40 **Summary**

41 The function of nearly half of all protein-coding genes identified in bacterial genomes remains  
42 unknown. To systematically explore the functions of these proteins, we generated saturated  
43 transposon mutant libraries from 25 diverse bacteria and we assayed mutant phenotypes  
44 across hundreds of distinct conditions. From 3,903 genome-wide mutant fitness assays, we  
45 obtained 14.9 million gene phenotype measurements and we identified a mutant phenotype for  
46 8,487 proteins with previously unknown functions. The majority of these hypothetical proteins  
47 (57%) had phenotypes that were either specific to a few conditions or were similar to that of  
48 another gene, thus enabling us to make informed predictions of protein function. For 1,914 of  
49 these hypothetical proteins, the functional associations are conserved across related proteins  
50 from different bacteria, which confirms that these associations are genuine. This comprehensive  
51 catalogue of experimentally-annotated protein functions also enables the targeted exploration of  
52 specific biological processes. For example, sensitivity to a DNA-damaging agent revealed 28  
53 known families of DNA repair proteins and 11 putative novel families. Across all sequenced  
54 bacteria, 14% of proteins that lack detailed annotations have an ortholog with a functional  
55 association in our data set. Our study demonstrates the utility and scalability of high-throughput  
56 genetics for large-scale annotation of bacterial proteins and provides a vast compendium of  
57 experimentally-determined protein functions across diverse bacteria.

58

## 59 **Background**

60 Tens of thousands of bacterial genomes have been sequenced, revealing the predicted amino  
61 acid sequences of millions of distinct proteins. In sharp contrast, only a small proportion of these  
62 proteins have been studied experimentally and for most, presumptive functions can only be  
63 predicted via their similarity to experimentally characterized proteins. However, about one third  
64 of bacterial proteins are not similar enough to any characterized protein to be annotated by this  
65 approach. Furthermore, these predictions are often incorrect as homologous proteins may have  
66 different substrate specificities<sup>1</sup>. This sequence-to-function gap represents a growing challenge  
67 for microbiology, because new bacterial genomes are being sequenced at an ever-increasing  
68 rate, while experimental protein characterization continues to be relatively slow<sup>2</sup>.

69

70 One approach for determining an unknown protein's function is to assess the consequences of  
71 a loss-of-function mutation of the corresponding gene under a large number of conditions<sup>3-6</sup>.  
72 Transposon mutagenesis followed by sequencing (TnSeq) makes it possible to apply this  
73 principle more systematically and to identify mutant phenotypes at a genome-wide scale. In this  
74 approach, a transposon is used to create a complex library of mutant strains, each with a  
75 different random insertion site in the genome. If an insertion lies within a protein-coding gene, it  
76 disrupts the respective protein's function. DNA sequencing can identify where each transposon  
77 has inserted and quantify the abundance of each mutant strain. Importantly, pools of hundreds  
78 of thousands of different mutants can be grown together, enabling genome-wide measurements  
79 of phenotypes and protein function in a single experiment<sup>7,8</sup>. Coupling the approach to random  
80 DNA barcoding (RB-TnSeq) of individual mutants further increases the efficiency, facilitating the  
81 evaluation of a given mutant library under a large number of experimental conditions<sup>9</sup>. Here, we  
82 use RB-TnSeq to directly address the sequence-to-function gap by systematically exploring the  
83 phenotypes of thousands of proteins from 25 bacteria under hundreds of experimental  
84 conditions (**Fig. 1a**).

85

## 86 Results

### 87 **Mutant fitness compendia for 25 bacteria**

88

89 To perform a systematic assessment of protein functions across a phylogenetically diverse set  
90 of bacteria, we chose 25 genetically tractable bacteria representing five different bacterial  
91 divisions and 16 different genera. In addition to the model bacterium *Escherichia coli* and 22  
92 other aerobic heterotrophs, we studied a strictly anaerobic sulfate-reducing bacterium  
93 (*Desulfovibrio vulgaris*) and a strictly photosynthetic cyanobacterium (*Synechococcus*  
94 *elongatus*) (**Fig. 1b**). We generated a randomly barcoded transposon mutant library in each of  
95 the 25 bacteria, seven of which were previously described<sup>9-11</sup>. For each mutant population, we  
96 used TnSeq to generate genome-wide maps of transposon insertion locations (Methods).  
97 Mutant populations were highly complex with 28,252 to 504,158 uniquely barcoded transposon  
98 insertions per genome, corresponding to 5 to 66 insertions for the typical (median) protein-  
99 coding gene. Genes that have very few or no transposon insertions are likely to be essential for  
100 viability in the conditions that were used to select the mutants. We identified 289 to 614  
101 essential proteins per bacterium (6-23% of the genes in each), with 10,766 essential proteins  
102 total (Methods; **Supplementary Table 1**; **Supplementary Note 1**). To assess how many of  
103 these essential proteins are poorly annotated, we classified all existing protein annotations for  
104 each genome as “Detailed TIGR role” (for subfamilies that are annotated with functional roles by  
105 TIGRFAMs), “Other detailed” (for other functional annotations), “Vague” (for annotations like  
106 “transporter”), or “Hypothetical” (for functionally uninformative annotations) (Methods,  
107 **Supplementary Fig. 1**). Across the 25 bacteria, we identified 667 hypothetical proteins and 480  
108 proteins with vague annotations that were essential (**Supplementary Table 1**).

109

110 To identify conditions suitable for mutant fitness profiling, we tested the growth of the wild-type  
111 bacteria in a wide range of conditions, including the utilization of 94 different carbon sources  
112 and 45 different nitrogen sources, and their inhibition by 55 stress compounds including  
113 antibiotics, surfactants, and heavy metals (**Supplementary Tables 2-4**). If we identified  
114 utilization or inhibition, then we performed a genome-wide fitness experiment using the  
115 corresponding mutant library. In the typical experiment, we grew a pool of mutants for 4-8  
116 generations and used DNA barcode sequencing (BarSeq)<sup>12</sup> to compare the abundance of the  
117 mutants before and after growth (**Fig. 1a**). For each gene, we defined gene fitness to be the log<sub>2</sub>  
118 change in abundance of mutants in that gene during the experiment (Methods, **Fig. 1a**).

119

120 For each bacterium, we successfully assayed fitness in 27 to 129 different experimental  
121 conditions (**Fig. 1b**), with a total of 2,035 bacterium-condition combinations. These conditions  
122 included the growth conditions discussed above as well as growth under varying pH and  
123 temperature and motility on agar plates. Including replicates, we conducted a total of 3,903  
124 genome-wide fitness experiments that met our criteria for biological and internal consistency<sup>9</sup>,  
125 and we obtained 14.9 million gene fitness values for 96,024 different non-essential protein-  
126 coding genes. The conditions and the fitness data can be viewed at the Fitness Browser  
127 (<http://fit.genomics.lbl.gov/>). Taken together, our results provide high-resolution genome-wide  
128 fitness maps for all 25 bacteria examined.

129

130 To establish the consistency of our data with known protein functions, we examined fitness data  
131 for the three most common classes of experiments: carbon utilization, nitrogen utilization, and  
132 stress. Genes without phenotypes have fitness values near 0, genes that are important for  
133 fitness have fitness values less than 0, and genes that are detrimental to fitness have fitness  
134 values greater than 0 (**Fig. 1c**). For the utilization of D-fructose or 4-hydroxybenzoate as the  
135 sole source of carbon in *Cupriavidus basilensis*, the fitness data identifies expected proteins for  
136 the catabolism of each substrate (**Fig. 1c**). Similarly, key enzymes and transporters required for  
137 the utilization of D-alanine or cytosine as the sole nitrogen source in *Azospirillum brasilense* are  
138 identified (**Supplementary Fig. 2a**). Lastly, in *Shewanella loihica*, orthologs of the CzcCBA

139 heavy metal efflux pump<sup>13</sup> and the zinc responsive regulator ZntR are important for fitness in  
140 the presence of elevated zinc (**Supplementary Fig. 2b**). In addition to identifying expected  
141 proteins, in all of the presented examples we also identified proteins that were previously not  
142 known to be involved in the respective process, including an efflux pump important for 4-  
143 hydroxybenzoate utilization by *C. basilensis* (**Fig. 1c, Supplementary Table 5**). In summary,  
144 these examples highlight how the fitness data can validate protein annotations and identify new  
145 proteins involved in diverse biological processes.

146  
147 We identified a significant mutant phenotype (false discovery rate < 5%) in at least one  
148 condition for 28,674 of the 96,024 (30%) non-essential proteins that we collected data for (**Fig.**  
149 **1b**). Proteins with high or moderate similarity to another protein in the same genome (paralogs,  
150 alignment score above 30% of the self-alignment score) are less likely to have a phenotype  
151 (25% vs. 31%,  $P < 10^{-15}$ ; **Supplementary Fig. 3**), which likely reflects genetic redundancy. 18%  
152 of all proteins with fitness measurements were detrimental to fitness in at least one condition  
153 (**Supplementary Fig. 3**), which is consistent with previous reports that many proteins are  
154 detrimental in some conditions<sup>3,14</sup>. Proteins annotated with a detailed TIGR role were  
155 particularly likely to have phenotypes, with more than half (52%) of those with fitness data  
156 showing a significant phenotype (**Fig. 1d**). In contrast, proteins that are not annotated with a  
157 detailed function (vague or hypothetical annotations) were less likely to have phenotypes (27%  
158 or 19%, respectively). Nevertheless, our assays identified phenotypes for 4,585 hypothetical  
159 proteins and for 3,902 other proteins with vague annotations (**Fig. 1d**). These include 2,828  
160 proteins that do not belong to any characterized family in either Pfam or TIGRFAMs<sup>15,16</sup>.  
161 Overall, we identified mutant phenotypes for 8,487 proteins that are not currently annotated with  
162 a detailed function, highlighting that a substantial proportion of hitherto vaguely annotated or  
163 hypothetical proteins have critical functions in bacteria.

### 164 165 **Conserved phenotypes are accurate predictors of protein function**

166  
167 To illuminate the biological function of individual proteins using genome-wide fitness data from  
168 multiple bacteria across a large number of experimental conditions, we used two principal  
169 approaches: (1) Identification of “specific” phenotypes that are observed only under one or a  
170 small number of conditions; (2) “Cofitness” patterns, where multiple proteins show similar fitness  
171 profiles across all conditions. Furthermore, we identified conserved specific phenotypes and  
172 conserved cofitness by comparing the data from 25 bacteria; and we tested the reliability of  
173 these conserved associations for annotating protein function.

174  
175 To assign specific phenotypes, we identified proteins that had a significant and notable  
176 phenotype ( $|\text{fitness}| > 1$ ) under only one or very few conditions (Methods). For example, the  
177 fluoride efflux protein CrcB<sup>17</sup> is required for fitness under elevated fluoride stress in multiple  
178 bacteria, but not for fitness in any of the hundreds of other experimental conditions that we  
179 tested (**Fig. 2a**). Among all genes with a significant phenotype under any condition, 34% have a  
180 specific phenotype (9,905 proteins), while the remaining genes either have weak but significant  
181 phenotypes ( $|\text{fitness}| < 1$ ) or complex phenotypic patterns.

182  
183 To assess whether specific phenotypes are useful indicators of protein function, we used our  
184 understanding of *E. coli* physiology and asked whether specific phenotypes in this bacterium  
185 could be used to accurately assign proteins to a “direct” enzymatic, regulatory, or transport  
186 function in a catabolic pathway. For example, all 4 proteins in *E. coli* with a specific phenotype  
187 during growth on D-xylose are directly involved in its utilization (*xyLABFR*). More broadly, we  
188 manually examined 101 experimentally characterized proteins according to the EcoCyc  
189 database<sup>18</sup> with specific phenotypes on individual carbon and nitrogen sources and found that  
190 67% have a known direct function in either the uptake or catabolism of that compound, or in the  
191 activation of genes that are (**Supplementary Table 6**). We also expected that if the specific  
192 phenotype were conserved (an orthologous protein from another bacterium has a notable

193 phenotype in the same condition, see Methods), the association would be more reliable. Indeed,  
194 among the same 101 *E. coli* proteins, we found that 87% of proteins with conserved specific  
195 phenotypes were directly involved in utilization, as compared to 35% of proteins with specific  
196 phenotypes that were not conserved (**Fig. 2b**;  $P < 10^{-4}$ , Fisher exact test; **Supplementary**  
197 **Table 6**). Thus, if a protein has a specific phenotype in a defined medium, it is likely to be  
198 directly involved in utilizing that substrate and hence a protein function can be assigned, and the  
199 probability is higher if the phenotype is conserved in another bacterium.

200  
201 Overall, we identified specific phenotypes and conserved specific phenotypes for proteins of all  
202 annotation classes (**Fig. 2c**). In particular, specific phenotypes linked 2,970 proteins with vague  
203 or hypothetical annotations to over 100 diverse conditions, including 79 carbon sources, 42  
204 nitrogen sources, and 54 stresses. These include conserved specific phenotypes, and hence  
205 high-confidence associations, for 733 such poorly annotated proteins.

206  
207 Our second strategy for inferring a protein's role was based on the observation that proteins  
208 with related functions often have similar fitness patterns across multiple conditions, which we  
209 term "cofitness" <sup>3,4</sup>. For example, Npr and PtsP of the nitrogen phosphotransferase system  
210 (PTS) in *E. coli* exhibit high cofitness (**Fig. 2d**), consistent with the known role of PtsP in the  
211 phosphorylation and activation of Npr <sup>19</sup>. Taking advantage of our comprehensive fitness  
212 dataset across 25 bacteria, we can now systematically address whether conserved cofitness is  
213 a stronger indicator of function than cofitness in one bacterium. For example, the orthologs of  
214 Npr and PtsP also have high cofitness in *S. oneidensis* (**Fig. 2d**).

215  
216 We first identified all pairs of proteins with highly correlated phenotypes within an organism (cofit  
217 proteins), and the subset of these with orthologous protein pairs that are cofit across two or  
218 more organisms (conserved cofit proteins). To test how accurately cofitness or conserved  
219 cofitness links together functionally-related proteins, we next identified protein pairs with existing  
220 detailed functional annotations (a TIGR subrole from the TIGRFAMs database of protein  
221 families <sup>16</sup>) and determined the frequency with which the subrole annotation of a protein is  
222 accurately "predicted" by that of the highest scoring cofit protein that is not nearby (see Methods  
223 for details). High-scoring cofitness in a single bacterium leads to predictions of TIGR subroles  
224 that are mostly correct, but the accuracy decays rapidly as the cofitness score decreases (**Fig.**  
225 **2e**). In contrast, for conserved cofitness, the decay is much slower (**Fig. 2e**). Furthermore,  
226 conserved cofitness is significantly more accurate for a given number of predictions: for  
227 example, the top 2,000 predictions from cofitness ( $r > 0.81$  for gene pairs from one bacterium)  
228 have 62% agreement, while the top 2,000 predictions from conserved cofitness ( $r > 0.56$  for  
229 gene pairs from both bacteria) have 69% agreement ( $P < 10^{-5}$ , Fisher exact test). Using  
230 thresholds of  $r > 0.8$  for cofitness or  $r > 0.6$  for conserved cofitness, we identified at least one  
231 association from cofitness or conserved cofitness for 15% of the genes with fitness data and for  
232 45% of genes with significant phenotypes. We identified associations for all types of proteins  
233 (**Fig. 2f**), including for 1,934 hypothetical proteins and for 1,674 other vaguely-annotated  
234 proteins. When the associations link proteins that lack detailed annotations to proteins with  
235 TIGR subroles, the top-level roles are diverse, with the most common role ("cellular processes")  
236 accounting for just 21% of them.

237  
238 Combining the specific phenotypes and cofitness-derived functional associations, we identified  
239 a functional relationship for 19,962 proteins including 5,512 proteins with vague or hypothetical  
240 annotations. For 1,914 of these poorly-annotated proteins, we identified conserved associations,  
241 which are expected to be particularly robust predictors of protein function (**Supplementary**  
242 **Table 7**). These results demonstrate the utility of specific phenotypes and cofitness for  
243 elucidating the diverse roles of thousands of poorly-annotated proteins and highlight the  
244 advantages of collecting genome-wide fitness data for multiple bacteria.

245  
246 **Genetic overviews of biological processes across diverse bacteria**



247  
248 Genome-wide mutant fitness profiling of phylogenetically diverse bacteria across a wide range  
249 of biological conditions provides a comprehensive genetic overview of each biological condition  
250 studied. To illustrate this, we examined proteins with conserved, specific, and important  
251 phenotypes (fitness < 0) during cisplatin stress (**Fig. 3a**). Cisplatin reacts with DNA to form  
252 crosslinks that block DNA replication, so we expected that DNA repair proteins would be  
253 important for growth in this condition<sup>3</sup>. Indeed, of the 57 protein families that were specifically  
254 important for resisting cisplatin in more than one bacterium, 28 are known to be involved in DNA  
255 repair including the UvrABC nucleotide excision repair complex and the RecFOR homologous  
256 recombination pathway (**Fig. 3a, Supplementary Table 8**). In addition, 6 of the other  
257 characterized families that have conserved sensitivity to cisplatin are involved in cell division or  
258 chromosome segregation (**Fig. 3a**), which is relevant because DNA damage can inhibit DNA  
259 replication and lead to filamentous cells<sup>20</sup>.

260  
261 We predict that 11 of the remaining 23 families with conserved cisplatin phenotypes are also  
262 involved in DNA repair, either because they contain DNA-related domains or because similar  
263 proteins are regulated by the DNA damage response regulator LexA in some bacteria<sup>21-23</sup>.  
264 Three of the putative novel DNA repair genes (*endA*, *yejH*, *yhgF*) are present in the well-studied  
265 bacterium *E. coli*, and all 3 are important for resistance to ionizing radiation, which also  
266 damages DNA<sup>24</sup>. This strongly suggests that their phenotypes are due to DNA damage. We  
267 also confirmed that an *E. coli* strain that has *endA* deleted is sensitive to cisplatin  
268 (**Supplementary Fig. 4**). The remaining 12 families likely have indirect roles in DNA repair,  
269 including 3 genes involved in tRNA or rRNA metabolism. Overall, our cisplatin experiments  
270 provide an overview of the proteins involved in DNA repair across diverse bacteria, including 28  
271 protein families with a known role and 11 protein families whose putative role in DNA repair is  
272 not yet understood.

273  
274 Similar genetic overviews were systematically obtained for the wide range of metabolic  
275 processes studied. To illustrate this, we identified at least one protein with a conserved, specific,  
276 and important phenotype in 67 carbon sources and 39 nitrogen sources. As an illustrative  
277 example, we examined D-xylose catabolism, which we assayed as the sole carbon source in 8  
278 bacteria. We found that XylAB is required in *E. coli* and in 6 other bacteria, confirming its central  
279 and conserved role in D-xylose catabolism (**Fig. 3b**). In contrast, the well-characterized *E. coli*  
280 XylR regulator and XylF transporter are not required in each of the other 6 bacteria: two  
281 *Pseudomonads* use alternative transport proteins for D-xylose while *Phaeobacter inhibens* and  
282 *Sinorhizobium meliloti* require a *lacI*-like regulator for D-xylose utilization, as previously  
283 predicted<sup>9,25</sup>. In contrast to the 7 bacteria that use the canonical XylAB pathway, we found that  
284 *Sphingomonas koreensis* uses an alternative, oxidative pathway for D-xylose utilization<sup>26,27</sup>.  
285 Fitness data and comparative analysis to a similar pathway in archaea<sup>26,28</sup> suggest that *xylX*  
286 encodes the required enzyme 2-keto-3-deoxyxylonate dehydratase in *S. koreensis*  
287 (**Supplementary Table 9**). Our analysis also identified additional genes including *lon*, *galM*, and  
288 an *icIR*-like regulator that show conserved D-xylose utilization phenotypes across multiple  
289 bacteria, but whose exact roles in D-xylose catabolism remain to be elucidated. This  
290 comparative phenotypic analysis of D-xylose as a single carbon source represents just one of  
291 more than 100 carbon and nitrogen source experiments studied and highlights the power of our  
292 approach for the validation and discovery of new proteins and pathways involved in basic  
293 metabolic processes.

#### 294 295 **Accurate annotations of individual protein families**

296  
297 Large-scale mutant fitness data can also be used to improve our understanding of proteins  
298 annotated with a general biochemical function but lacking substrate specificity. To illustrate this,  
299 we used the mutant fitness data to systematically reannotate the substrate specificities of 101  
300 ABC transporter family proteins that have strong and specific phenotypes (fitness < -2 and

301 statistically significant) during the utilization of diverse carbon or nitrogen sources (**Fig. 3c**;  
302 **Supplementary Table 10**). 20 of the 101 proteins only have vague annotations, and for these  
303 we made novel substrate predictions based on the specific phenotype data. For example,  
304 Dshi\_0548 and Dshi\_0549 from *Dinoroseobacter shibae* are annotated with no substrate  
305 specificity yet are important for utilizing xylitol. For another 31 proteins with moderately-specific  
306 substrate annotations, such as "amino acids", we predicted a specific substrate within that  
307 group of compounds. For example, Ac3H11\_2942 and Ac3H11\_2943 from *Acidovorax* sp. 3H11  
308 are annotated as transporting "various polyols", whereas our data shows they are important for  
309 utilizing the polyol D-sorbitol but not the polyol D-mannitol. Another 14 proteins had incorrect  
310 annotations: for example, in three *Pseudomonads*, the L-carnitine transporter is mis-annotated  
311 as a choline transporter. In 10 cases, the specific phenotypes indicate that the protein transports  
312 a substrate that was not included in the annotation along with a substrate that was expected.  
313 For example, PS417\_12705 from *P. simiae* is annotated as transporting D-mannitol, but this  
314 protein is also important for utilizing D-sorbitol and D-mannose. 24 proteins had correct  
315 annotations that were confirmed by our data; these included 22 cases in which the specific  
316 phenotype(s) were expected given the annotation, and 2 cases in which the gene has the  
317 expected mutant phenotype(s) but the association to a specific condition was misleading. The  
318 data for the remaining 2 proteins was hard to interpret. Overall, our fitness data provided  
319 improved annotations for 75 of the 101 transporters examined. This analysis also highlights how  
320 detailed computational annotations are often misleading: for 24 of the 50 ABC transporters that  
321 were annotated with a substrate (48%), the predicted specificity was erroneous or incomplete.

### 322 323 **Annotation of uncharacterized protein families**

324  
325 To evaluate the utility of fitness data for elucidating the functions of entirely uncharacterized  
326 protein families, we identified conserved cofitness or a conserved specific phenotype for 288  
327 proteins that represent 77 different domains of unknown function (DUFs) from the Pfam  
328 database<sup>15</sup> (**Supplementary Table 11**). We manually examined the phenotypes of these 77  
329 DUFs and we propose broad functional annotations for 13 of them and specific molecular  
330 functions for an additional 8 (**Supplementary Note 2**). For example, proteins containing  
331 UPF0126 are specifically important for glycine utilization in five bacteria (**Fig. 4a**). Since this  
332 family is predicted to be a membrane protein, we propose that it is more specifically a glycine  
333 transporter. As a second example, we found that members of the UPF0060 family of predicted  
334 transmembrane proteins have a conserved specific phenotype in the presence of elevated  
335 thallium in three bacteria (**Fig. 4b**). Consequently, we propose that UPF0060-containing  
336 proteins may function as a thallium-specific efflux pump. As a third example, we found that in  
337 three bacteria, DUF2849-containing proteins are cofit with an adjacent sulfite reductase (*cysI*)  
338 gene (**Fig. 4c**). The sulfite reductase is important for fitness in most defined media conditions  
339 because it is involved in sulfate assimilation, which was the only source of sulfur in our base  
340 media, and DUF2849 also seems to be involved in this process. These three bacteria lack *cysJ*,  
341 which is usually the electron source for *cysI*, and other bacterial genomes that contain DUF2849  
342 also contain *cysI* but not *cysJ*. Based on this evidence, we propose that DUF2849 is an  
343 alternate electron source for sulfite reductase.

344  
345 Conserved phenotypic association can also be used to identify more complex pathways  
346 containing uncharacterized proteins. For example, the uncharacterized proteins YeaH and YcgB  
347 and the poorly characterized protein kinase YeaG<sup>29</sup> had high cofitness across six different  
348 bacteria (**Fig. 4d** shows the data from three of these bacteria; all  $r > 0.7$ ). Interestingly, while  
349 these three proteins are more cofit with each other than with any other protein in each of the six  
350 bacteria, the phenotypes are not conserved across species (**Fig. 4d**). We validated some of  
351 these key phenotypes in growth assays with individual mutants. These experiments confirmed  
352 that in *E. coli*, mutants in all three genes have a growth advantage when L-arginine is the  
353 nitrogen source (**Supplementary Fig. 5**), whereas in *S. oneidensis*, mutants in all three genes  
354 grow slowly when N-acetylglucosamine is the carbon source (**Supplementary Fig. 6**). Based on

355 these data and the protein kinase activity of YeaG, we propose that these three proteins act  
356 together in a conserved signaling pathway that is required for distinct cellular functions in  
357 different bacteria. Taken together, our data highlight how comprehensive fitness data can be  
358 used to provide novel experimental annotations for uncharacterized protein families and  
359 pathways.

360

### 361 **Relevance to all bacteria**

362

363 Beyond the 25 bacteria experimentally studied in the present study, combining fitness data with  
364 comparative genomics offers the opportunity to assign phenotype-derived protein annotations  
365 across all sequenced bacterial genomes. To address this, we focused on the utility of our  
366 dataset to illuminate the functions of hypothetical or vaguely-annotated proteins from  
367 phylogenetically diverse bacterial genomes (**Methods**). For poorly-annotated proteins from  
368 these bacterial genomes, 14% have a putative ortholog (with at least 30% sequence similarity)  
369 with a significant phenotype in our dataset (**Fig. 5a**). Furthermore, for 11% of the poorly-  
370 annotated proteins, we can associate an ortholog (above 30% similarity) to a specific condition  
371 or to another protein with cofitness (**Fig. 5a**). Even at 30% similarity, our data should provide  
372 functional insight for many poorly characterized bacterial proteins. Supporting this, using gene  
373 ontology annotations supported by experimental evidence, Clark and Radivojac analyzed the  
374 similarity of molecular function between homologous proteins and found 60-70% functional  
375 conservation for homologs from anywhere from 30-100% identity<sup>30</sup>. Not surprisingly, the  
376 probability of finding an ortholog with a functional association in our dataset is much higher for  
377 poorly-annotated proteins from bacterial divisions that we studied multiple representatives of  
378 ( $\alpha$ , $\beta$ , $\gamma$ -Proteobacteria) than for such proteins from other bacteria (21% vs. 6%). For bacteria  
379 from these three divisions, we provide putative orthologs with functional associations for about  
380 330 poorly-annotated proteins per average genome (3,913 proteins per genome \* 40% poorly-  
381 annotated \* 21%). To enable rapid access to the phenotypes and functional associations of the  
382 homologs of a protein of interest, we provide a "Fitness BLAST" web service  
383 ([http://fit.genomics.lbl.gov/images/fitblast\\_example.html](http://fit.genomics.lbl.gov/images/fitblast_example.html)). The results from Fitness BLAST are  
384 available within the protein pages at IMG/M<sup>31</sup> and MicrobesOnline<sup>32</sup>. We also provide a web  
385 page for comparing all of the proteins in a bacterium to the fitness data. Overall, our data  
386 provides functional information by homology for 14% of poorly annotated proteins across all  
387 sequenced bacteria.

388

389

### 390 **Discussion**

391

392 We have shown that high-throughput genetics can provide mutant phenotypes and functional  
393 associations for thousands of vaguely annotated and hypothetical bacterial proteins. Many of  
394 these functional associations were conserved, which increases the reliability of these  
395 associations. These associations can also highlight errors in current computational annotations,  
396 as we demonstrated for ABC transporters. Additionally, our comparative data provides unbiased  
397 functional overviews of biological processes by identifying proteins that are important for fitness  
398 under the same condition in multiple bacteria, as illustrated for D-xylose utilization and cisplatin  
399 stress.

400

401 The major challenge in extending our results to all bacterial proteins is their incredible diversity.  
402 Although we identified functional associations for 19,962 proteins and for 5,512 proteins that  
403 lack detailed annotations, these include putative orthologs of just 11% of the bacterial proteins  
404 that lack detailed annotations. Improving this coverage will require a larger effort to generate  
405 mutants in more diverse bacteria: our study included representatives of only 5 of the ~40  
406 divisions of bacteria that have been cultivated so far. Another challenge will be to improve the  
407 inference of protein function from phenotype. Many of the vaguely-annotated or incorrectly-  
408 annotated proteins belong to characterized families of enzymes or transporters and the main



409 uncertainty as to their function is what substrate they act on. We proposed functional  
410 associations for 64% of the proteins that have significant phenotypes, and in principle, these  
411 associations could be used to automatically identify the physiological substrate for many of  
412 these proteins.

413  
414 For thousands of proteins that previously lacked an informative annotation, our mutant  
415 phenotypes, and the functional associations derived from them, provide a rich resource to guide  
416 further study. To facilitate this, we developed the Fitness Browser web site  
417 (<http://fit.genomics.lbl.gov>) to view the fitness data for a gene or condition of interest. This site  
418 also supports the comparison of the fitness data across bacteria and incorporates tools for  
419 sequence-based annotation. In summary, our study demonstrates the scale with which large-  
420 scale fitness data can be collected in diverse bacteria and the utility of these data to provide  
421 insights into the functions of thousands of proteins, thereby helping to close the sequence-to-  
422 function gap in microbiology.

423  
424

## 425 **Methods**

426 **Bacteria.** The bacteria mutagenized in this study are listed in **Supplementary Table 12**. Seven  
427 bacteria were isolated from groundwater collected from different monitoring wells at the Oak  
428 Ridge National Laboratory Field Research Center (FRC; <http://www.esd.ornl.gov/orifrc/>), and  
429 five have not been described previously: *Acidovorax* sp. GW101-3H11, *Pseudomonas*  
430 *fluorescens* FW300-N1B4, *P. fluorescens* FW300-N2E3, *P. fluorescens* FW300-N2C3, and *P.*  
431 *fluorescens* GW456-L13. *Acidovorax* sp. GW101-3H11 was isolated as a single colony on a  
432 Luria-Bertani (LB) agar plate grown at 30°C using an inoculum from FRC well GW101. *P.*  
433 *fluorescens* FW300-N1B4, *P. fluorescens* FW300-N2E3, and *P. fluorescens* FW300-N2C3 were  
434 all isolated at 30°C under anaerobic denitrifying conditions with acetate, propionate, and  
435 butyrate as the carbon source, respectively, using inoculum from FRC well FW300.  
436 *Pseudomonas fluorescens* GW456-L13 was isolated from FRC well FW456 under anaerobic  
437 incubations on a LB agar plate. We previously described the isolation of *Pseudomonas*  
438 *fluorescens* FW300-N2E2<sup>33</sup> and *Cupriavidus basilensis* 4G11<sup>34</sup>. We also studied individual  
439 mutants of several organisms. For *E. coli* BW25113, we used single-gene deletions from the  
440 Keio collection<sup>35</sup>. For *S. oneidensis* MR-1, we used transposon mutants that had been  
441 individually sequenced<sup>4</sup>.

442  
443 **Media and standard culturing conditions.** A full list of the medias used in this study and their  
444 components are given in **Supplementary Table 13**. We routinely cultured *Acidovorax* sp.  
445 GW101-3H11, *Azospirillum brasilense* Sp245, *Burkholderia phytofirmans* PsJN, *Escherichia coli*  
446 BW25113, all *Pseudomonads* and *Shewanellae*, *Sinorhizobium meliloti* 1021, and  
447 *Sphingomonas koreensis* DSMZ 15582 in LB. *Cupriavidus basilensis* 4G11 was typically  
448 cultured in R2A media. *Dechlorosoma suillum* PS was cultured in ALP media<sup>11</sup>. *Desulfovibrio*  
449 *vulgaris* Miyazaki F was grown anaerobically in lactate-sulfate (MOLS4) media, as previously  
450 described<sup>36,37</sup>. We used marine broth (Difco 2216) for standard culturing of *Dinoroseobacter*  
451 *shibae* DFL-12, *Kangiella aquimarina* SW-154T, *Marinobacter adhaerens* HP15, and  
452 *Phaeobacter inhibens* BS107. *Synechococcus elongatus* PCC 7942 was normally cultured in  
453 BG-11 media with either 7,000 or 9,250 lux. All bacteria were typically cultured at 30°C except  
454 *Escherichia coli* BW25113 and *Shewanella amazonensis* SB2B, which were cultured at 37°C,  
455 and *P. inhibens* BS107, which was grown at 25°C. The *E. coli* conjugation strain WM3064 was  
456 cultured in LB at 37°C with diaminopimelic acid (DAP) added to a final concentration of 300 µM.

457  
458 **High-throughput growth assays of wild-type bacteria.** To assess the phenotypic capabilities  
459 of 23 aerobic heterotrophic bacteria and to identify conditions suitable for mutant fitness  
460 profiling, we monitored the growth of the wild-type bacterium in a 96-well microplate assay.

461 These prescreen growth assays were performed in a Tecan microplate reader (either Sunrise or  
462 Infinite F200) with absorbance readings (OD<sub>600</sub>) every 15 minutes. All 96-well microplate growth  
463 assays contained 150  $\mu$ L culture volume per well at a starting OD<sub>600</sub> of 0.02. We used the grofit  
464 package in R<sup>38</sup> to analyze all growth curve data in this study. For carbon and nitrogen source  
465 utilization, we tested 94 and 45 possible substrates, respectively, in a defined medium  
466 (**Supplementary Tables 2 and 3**). We classified a bacterium as positive for usage of a  
467 particular substrate if (1) the maximum OD<sub>600</sub> on the substrate was at least 1.5 greater than the  
468 average of the water controls and the integral under the curve (spline.integral) was 10% greater  
469 than the average of the water controls or (2) a successful genome-wide fitness assay was  
470 collected on the substrate, as described below. We included the second criterion because our  
471 automated scoring of the wild-type growth curves was conservative and did not include all  
472 conditions used for genome-wide fitness assays.

473  
474 Additionally, for each of the 23 heterotrophic bacteria, we determined the inhibitory  
475 concentrations for 39-55 diverse stress compounds including antibiotics, biocides, metals,  
476 furans, aldehydes, and oxyanions. For each compound, we grew the wild-type bacterium across  
477 a 1,000-fold range of inhibitor concentrations in a rich media. We used the spline.integral  
478 parameter of grofit to fit dose-response curves and calculate the half-maximum inhibitory  
479 concentration (IC<sub>50</sub>) values for each compound (**Supplementary Table 4**). For *Desulfovibrio*  
480 *vulgaris* Miyazaki F and *Synechococcus elongatus* PCC 7942, we did not perform these growth  
481 prescreen assays, rather, we just performed the mutant fitness assays across a broad range of  
482 inhibitor concentrations.

483  
484 **Genome sequencing.** We sequenced *Acidovorax sp.* GW101-3H11 and five *Pseudomonads*  
485 by using a combination of Illumina and Pacific Biosciences. For Illumina-first assembly, we used  
486 scythe (<https://github.com/vsbuffalo/scythe>) and sickle (<https://github.com/najoshi/sickle>) to trim  
487 and clean Illumina reads, we assembled with SPAdes 3.0<sup>39</sup>, we performed hybrid  
488 Illumina/PacBio assembly on SMRTportal using AHA<sup>40</sup>, we used BridgeMapper on SMRTportal  
489 to fix misassembled contigs, we mapped Illumina reads to the new assembly with bowtie 2<sup>41</sup>,  
490 and we used pilon to correct local errors<sup>42</sup>. For *Acidovorax sp.* GW101-3H11, we instead used  
491 A5<sup>43</sup> to assemble the Illumina reads and we used AHA to join contigs together. For PacBio-first  
492 assembly, we used HGAP3 on SMRTportal, we used circlator to find additional joins<sup>44</sup>, and we  
493 again used bowtie 2 and pilon to correct local errors. See **Supplementary Table 14** for a  
494 summary of these genome assemblies and their accession numbers. In addition,  
495 *Sphingomonas koreensis* DSMZ 15582 was sequenced for this project by the Joint Genome  
496 Institute, using Pacific Biosciences.

497  
498 **Constructing pools of randomly barcoded transposon mutants.** The transposon mutant  
499 libraries for seven bacteria were described previously<sup>9-11</sup>. The other 18 bacteria were  
500 mutagenized with randomly barcoded plasmids containing a *mariner* or *Tn5* transposon, a *pir*-  
501 dependent conditional origin of replication, and a kanamycin resistance marker, using the  
502 vectors that we described previously<sup>9</sup>. The plasmids were delivered by conjugation with *E. coli*  
503 WM3064, which is a diaminopimelate auxotroph and is *pir*<sup>+</sup>. The conditions for mutagenizing  
504 each organism are described in **Supplementary Table 15**. Generally, we conjugated mid-log  
505 phase grown WM3064 donor (either *mariner* donor plasmid library APA752 or *Tn5* donor  
506 plasmid library APA766) and recipient cells on 0.45  $\mu$ M nitrocellulose filters (Millipore) overlaid  
507 on rich media agar plates supplemented with DAP. We used the rich medium preferred by the  
508 recipient (**Supplementary Table 15**). After conjugation, filters were resuspended in recipient  
509 rich media and plated on recipient rich media agar plate supplemented with kanamycin. After  
510 growth, we scraped together kanamycin resistant colonies into recipient rich media with  
511 kanamycin, diluted the culture back to a starting OD<sub>600</sub> of 0.2 in 50-100 mL of recipient rich  
512 media with kanamycin, and grew the mutant library to a final OD<sub>600</sub> of between 1.0 and 2.0. We  
513 added glycerol to a final volume of 10%, made multiple 1 or 2 mL -80°C freezer stocks, and  
514 collected cell pellets to extract genomic DNA for TnSeq. For *Desulfovibrio vulgaris* Miyazaki F,

515 we selected for G418-resistant transposon mutants in liquid media with no plating step  
516 (**Supplementary Table 15**).

517  
518 **Transposon insertion site sequencing (TnSeq).** Given a pool of mutants, we performed  
519 TnSeq just once to amplify and sequence the transposon junction and to link the barcodes to a  
520 location in the genome<sup>9</sup>. We considered a barcode to be confidently mapped to a location if this  
521 mapping was supported by at least 10 reads. (For *Shewanella* sp. ANA-3, the threshold was 8  
522 reads.) The number of unique barcodes (strains) mapped in each mutant library is shown in  
523 **Supplementary Table 15**. Given this mapping, the abundance of the strains in each sample  
524 can be determined by a simpler and cheaper protocol: amplifying the barcodes with PCR  
525 followed by barcode sequencing<sup>12</sup>.

526  
527 **Identifying essential genes.** Genes that lack insertions or that have very low coverage in the  
528 start samples are likely to be essential or to be important for growth in rich media, as except for  
529 *Synechococcus elongatus*, pools of mutants were produced and recovered in media that  
530 contained yeast extract. We used previously published heuristics<sup>10</sup> to distinguish likely-essential  
531 genes from genes that are too short or that are too repetitive to map insertions in. Briefly, for  
532 each protein-coding gene, we computed the total read density in TnSeq (reads / nucleotides  
533 across the entire gene) and the density of insertion sites within the central 10-90% of each gene  
534 (sites/nucleotides). We then excluded genes that might be difficult to map insertions within  
535 because they were very similar to other parts of the genome (BLAT score above 50) and also  
536 very-short genes of less than 100 nucleotides. Given the median insertion density and the  
537 median length of the remaining genes, we asked how short a gene could be and still be unlikely  
538 to have no insertions at all by chance ( $P < 0.02$ , Poisson distribution). Genes shorter than this  
539 threshold were excluded; the threshold varied from 100 nucleotides for *Phaeobacter inhibens* to  
540 600 nucleotides for *Desulfovibrio vulgaris* Miyazaki. For the remaining genes, we normalized the  
541 read density by GC content by dividing by the running median of read density over a window of  
542 201 genes (sorted by GC content). We normalized the insertion density so that the median  
543 gene's value was 1. Protein-coding genes were considered essential or important for growth if  
544 we did not estimate fitness values for the gene and both the normalized insertion density and  
545 the normalized read density were under 0.2. A validation of this approach is described in  
546 **Supplementary Note 1**.

547  
548 **Mutant fitness assays.** For each mutant library, we performed competitive mutant fitness  
549 assays under a large number of growth conditions that were chosen based on the results of  
550 high-throughput growth assays of wild-type bacteria (see above). The full list of experiments  
551 performed for each mutant library is available at <http://fit.genomics.lbl.gov>. Our analysis includes  
552 385 successful experiments from Wetmore *et al.*<sup>9</sup> and 36 successful experiments from Melnyk  
553 *et al.*<sup>11</sup>. The other 3,482 successful fitness assays are described here for the first time. In  
554 general, all growth assays with carbon sources, nitrogen sources, and inhibitors were done as  
555 previously described<sup>9</sup>. Briefly, an aliquot of the mutant library was thawed and inoculated into  
556 25 mL of rich media with kanamycin and grown to mid-log phase in a flask. Depending on the  
557 mutant library, this growth recovery took between 3 and 24 hours. After recovery, we collected  
558 pellets for genomic DNA extraction and barcode sequencing (BarSeq) of the input or “start”  
559 sample. We used the remaining cells to set up multiple mutant fitness assays with diverse  
560 carbon and nitrogen sources in defined media and diverse inhibitors in rich media, all at a  
561 starting OD<sub>600</sub> of 0.02. In addition, for most bacteria, we profiled growth of the mutant library at  
562 different pH and at different temperatures. After the mutant library grew to saturation under the  
563 selective growth condition (typically 4 to 8 population doublings), we collected a cell pellet for  
564 genomic DNA extraction and BarSeq of the “end” sample. As described below, we calculate  
565 gene fitness from the barcode counts of the end sample relative to the start sample.

566  
567 We used a number of different growth formats and media formations for mutant fitness assays  
568 across the 25 bacteria. The complete metadata for all mutant fitness assays in each bacterium

569 are available at the supporting website. A full list of compound components for each growth  
570 media are contained in **Supplementary Table 13**. Many fitness assays were done in 48-well  
571 microplates (Greiner) with 700  $\mu$ L culture volume per well and grown in a Tecan Infinite F200  
572 plate reader with OD<sub>600</sub> measurements every 15 minutes. For these 48-well microplate assays,  
573 we combined the cultures from two replicate wells before genomic DNA extraction (total volume  
574 of experiment = 1.4 mL). For 24-well microplate experiments, we used deep-well plates with 1.5  
575 mL (for inhibitors) or 2 mL (for carbon and nitrogen sources) total culture volume per well. All  
576 24-well microplate experiments were grown in a Multitron incubating shaker (Innova). For 24-  
577 well microplate experiments, we typically took the OD<sub>600</sub> of each culture every 12 to 24 hours in  
578 a Tecan plate reader (after transferring the cells to a Greiner 96-well microplate). Over 1,000  
579 experiments, primarily carbon source and temperature experiments, were done in glass test  
580 tubes with 5 mL culture volumes. For the test tube experiments, we monitored OD<sub>600</sub> every 12 to  
581 24 hours with a standard spectrophotometer and cuvettes.

582  
583 For stress experiments, we tried to use a concentration of each compound that allows growth  
584 (because if there is no growth, then the abundance of the strains will not change) but  
585 significantly inhibits growth (or else the fitness pattern is likely to be as if the compound were not  
586 added). Ideally, the concentration is such that the growth rate is cut in half. For aerobic  
587 heterotrophs, we measured growth across several orders of magnitude of concentrations of  
588 each stress compound, as described above. This gave a rough estimate of what concentration  
589 to use. Then, when we performed the fitness assays, we used a few different concentrations to  
590 try and capture an inhibitory but sub-lethal concentration. For assays done in 48-well  
591 microplates and grown in a Tecan Infinite F200 plate reader, we could confirm that the culture  
592 was inhibited relative to a no stress control. For stress assays in 24-well, deep-well microplates  
593 and grown in the Multitron shaker, we took OD readings approximately every 12 hours to  
594 estimate which cultures were inhibited. In practice, we often did multiple mutant fitness assays  
595 with different concentrations of the same inhibitor. We also collected fitness data in plain rich  
596 media without an added inhibitory compound.

597 For some carbon source experiments in *Desulfovibrio vulgaris* Miyazaki F, which is strictly  
598 anaerobic, we grew the mutant pool in 18 x 150 mm hungate tubes with a butyl rubber stopper  
599 and an aluminum crimp seal (Chemglass Life Sciences, Vineland, NJ) with a culture volume of  
600 10 mL and a headspace of about 15 mL. For the remainder of the *Desulfovibrio vulgaris* fitness  
601 experiments, we grew the mutant pool in 24 well microplates inside of the anaerobic chamber.  
602 We used OD<sub>600</sub> measurements to determine which cultures were inhibited by varying  
603 concentrations of stress compounds. Similarly, for six of the other heterotrophs, we measured  
604 gene fitness during anaerobic growth. All anaerobic media was prepared within a Coy anaerobic  
605 chamber with an atmosphere of about 2% H<sub>2</sub>, 5% CO<sub>2</sub>, and 93% N<sub>2</sub>.

606 For *Synechococcus elongatus* PCC 7942, which is strictly photosynthetic, we recovered the  
607 library from the freezer in BG-11 media at a light level of 7000 lux and we conducted fitness  
608 assays at 9250 lux. We used OD<sub>750</sub> to measure the growth of *S. elongatus*. Most *S. elongatus*  
609 mutant fitness assays were done in the wells of a 12-well microplate (Falcon) with a 5 mL  
610 culture volume.

611 In addition to growth assays in liquid media, we successfully studied motility in 12 bacteria using  
612 a soft agar assay. For motility assays, the mutant pool was inoculated into the center of a 0.3%  
613 agar rich media plate and “outer” samples with motile cells were removed with a razor after 24-  
614 48 hours. In many instances, we also removed an “inner” sample of cells from near the point of  
615 inoculation. Not all bacteria we assayed were motile in this soft agar assay and others were  
616 motile but did not give mutant fitness results that passed our quality metrics.

617



618 In four bacteria, we also assayed survival. In these assays, a mutant pool was subjected to a  
619 stressful condition (either extended stationary phase or a low temperature of 4°C) for a defined  
620 period; then, to determine which strains are still viable, they were recovered in rich media for a  
621 few generations. After recovery in rich media, the cells were harvested for genomic DNA  
622 extraction and BarSeq.

623 **Barcode sequencing (BarSeq).** Genomic DNA extraction and barcode PCR were performed  
624 as described previously<sup>9</sup>. Most genomic DNA extractions were done in a 96-well format using a  
625 QIAcube HT liquid handling robot (QIAGEN). We used the 98°C BarSeq PCR protocol<sup>9</sup>, which  
626 is less sensitive to high GC content. In general, we multiplexed 48 samples per lane of Illumina  
627 HiSeq. For *E. coli*, we sequenced 96 samples per lane instead.

628  
629 **Computation of fitness values.** Fitness data was analyzed as previously described<sup>9</sup>. Briefly,  
630 the fitness value of each strain (an individual transposon mutant) is the normalized  $\log_2$ (strain  
631 barcode abundance at end of experiment/strain barcode abundance at start of experiment). The  
632 fitness value of each gene is the weighted average of the fitness of its strains; only strains that  
633 have sufficient start reads and lie within the central 10-90% of the gene are included. The  
634 median number of usable strains per gene in each bacterium is shown in **Fig. 1b**. The gene  
635 fitness values were then normalized to remove the effects of variation in genes' copy number:  
636 the median for each scaffold is set to zero, and for large scaffolds, the running median of the  
637 gene fitness values is subtracted. Also, for large scaffolds, the peak of the distribution of gene  
638 fitness values is set to be at zero. For example, a gene fitness value of -2 means that the strains  
639 with transposon mutant insertions in that gene were, on average, at 25% of their original  
640 abundance at the end of the experiment.

641  
642 Fitness experiments were deemed successful using the quality metrics that we described  
643 previously<sup>9</sup>. These metrics ensure that the typical gene has sufficient coverage, that the fitness  
644 values of independent insertions in the same gene are consistent, and that there is no GC bias  
645<sup>9</sup>. Experiments that did not meet these thresholds were excluded from our analyses. The  
646 remaining experiments show good agreement between biological replicates, with a median  
647 correlation of 0.88 for gene fitness values from defined media experiments. Stress experiments  
648 sometimes have little biological signal, as they are usually done in rich media and mutants of  
649 genes that are important for growth in rich media may be absent from the pools. Nevertheless,  
650 the median correlation between replicate stress experiments was 0.68.

651  
652 To estimate the reliability of the fitness value for a gene in a specific experiment, we use a *t*-like  
653 test statistic which is the gene's fitness divided by the standard error<sup>9</sup>. The standard error is the  
654 maximum of two estimates. The first estimate is based on the consistency of the fitness for the  
655 strains in that gene. The second estimate is based on the number of reads for the gene.

656  
657 Even mild phenotypes were quite consistent between replicate experiments if they were  
658 statistically significant. For example, if a gene had a mild but significant phenotype in one  
659 replicate ( $0.5 < |\text{fitness}| < 2$  and  $|t| > 4$ ), then the sign of the fitness value was the same in the  
660 other replicate 95.5% of the time. Because this comparison might be biased if the two replicates  
661 were compared to the same control sample, only replicates with independent controls were  
662 included.

663  
664 **Genes with statistically significant phenotypes.** We averaged fitness values from replicate  
665 experiments. We combined *t* scores across replicate experiments with two different approaches.  
666 If the replicates did not share a start sample and were entirely independent, then we used  $t_{\text{comb}}$   
667  $= \text{sum}(t) / \text{sqrt}(n)$ , where *n* is the number of replicates. But if the replicates used the same start  
668 sample then this metric would be biased. To correct for this, we assumed that the start and end  
669 samples have similar amounts of noise. This is conservative because we usually sequenced the

670 start samples with more than one PCR and with different multiplexing tags. Given this  
671 assumption and given that  $\text{variance}(A+B) = \text{variance}(A) + \text{variance}(B)$  if A and B are  
672 independent random variables, it is easy to show that the above estimate of  $t_{\text{comb}}$  needs to be  
673 decreased by a factor of  $\sqrt{(n^2 + n)/(2n)}$ , where  $n$  is the number of replicates.

674  
675 Our standard threshold for a significant phenotype was  $|\text{fitness}| > 0.5$  and  $|t| > 4$ , but  
676 this was increased for some bacteria to maintain a false discovery rate (FDR) of less than 5%.  
677 We use a minimum threshold on  $|\text{fitness}|$  as well as  $|t|$  to account for imperfect normalization or  
678 for other small biases in the fitness values. To estimate the number of false positives, we used  
679 control experiments, that is, comparisons between different measurements of different aliquots  
680 of the same start sample. However we did not use some previously-published control  
681 comparisons (from <sup>9</sup>) that used the old PCR settings and had strong GC bias (to exclude these,  
682 we used the same thresholds that we used to remove biased experiments). The estimated  
683 number of false positive genes was then the number of control measurements that exceeded  
684 the thresholds, multiplied by the number of conditions and divided by the number of control  
685 experiments. As a second approach to estimate the number of false positives, we used the  
686 number of expected false positives if the  $t_{\text{comb}}$  values follow the standard normal distribution ( $2 * P(z > t) * \#\text{experiments} * \#\text{genes}$ ). If either estimate of the false discovery rate was above 5%,  
687 we raised our thresholds for both  $|\text{fitness}|$  and  $|t_{\text{comb}}|$  in steps of 0.1 and 0.5, respectively, until  
688 FDR < 5%. The highest thresholds used were  $|\text{fitness}| > 0.9$  and  $|t| > 6$ . Also, for *Pseudomonas*  
690 *fluorescens* FW300-N1B4, we identified six genes with large differences between control  
691 samples ( $|\text{fitness}| \approx 2$  and  $|t| \approx 6$ ). These genes cluster on the chromosome in two groups and  
692 are strongly co-fit, and several of the genes are annotated as being involved in capsular  
693 polysaccharide synthesis. Because this bacterium is rather sticky, we suspect that mutants in  
694 these genes are less adherent and were enriched in some control samples due to insufficient  
695 vortexing, so these six genes were excluded when estimating the number of false positives.

696  
697 **Sequence analysis.** To assign genes to Pfams <sup>15</sup> or TIGRFAMs <sup>16</sup>, we used HMMer 3.1b1 <sup>45</sup>  
698 and the trusted score cutoff for each family. We used Pfam 28.0 and TIGRFAM 15.0. We used  
699 only the curated families in Pfam ("Pfam A").

700  
701 To identify putative orthologs between pairs of genomes, we used bidirectional best protein  
702 BLAST hits with at least 80% alignment coverage both ways. We did not use any cutoff on  
703 similarity, as a similarity of phenotype can show that distant homologs have conserved  
704 functions.

705  
706 To measure the relevance of our fitness data to diverse bacteria, we started with 1,752 bacterial  
707 genomes from MicrobesOnline <sup>32</sup>. We grouped together closely related genomes if they were  
708 separated from their common ancestor by less than 0.01 substitutions per site in highly  
709 conserved proteins (using the MicrobesOnline species tree). These groups correspond roughly  
710 to species. We selected one representative of each group at random, which gave us 1,236  
711 divergent bacterial genomes. We selected 5 proteins at random, regardless of their annotation,  
712 from each of these genomes to form our sample of bacterial protein-coding genes. The best-  
713 scoring hit to one of the 25 bacteria that was studied was considered as a potential ortholog if  
714 the alignment coverage was 75% or more; we used a range of threshold for similarity but our  
715 recommended cutoff is that the ratio of the BLAST alignment score to the self score be above  
716 0.3 (as in <sup>46</sup>).

717  
718 To estimate the evolutionary relationships of the bacteria that we studied (**Fig. 1b**), we used  
719 Amphora2 <sup>47</sup> to identify 31 highly-conserved proteins in each genome and to align them, we  
720 concatenated the 31 protein alignments, and we used FastTree 2.1.8 <sup>48</sup> to infer a tree.

721  
722 **Specific phenotypes.** We defined a specific phenotype for a gene in an experiment as:  $|\text{fitness}|$   
723  $> 1$  and  $|t| > 5$  in this experiment;  $|\text{fitness}| < 1$  in at least 95% of experiments; and the fitness

724 value in this experiment is noticeably more extreme than most of its other fitness values  
725 ( $|fitness| > 95th\ percentile(|fitness|) + 0.5$ ).

726  
727 We considered a specific phenotype to be conserved if a potential ortholog had a specific  
728 phenotype with the same sign in a similar experiment with the same carbon source, nitrogen  
729 source, or stressful compound (but not necessarily using the same base media or the same  
730 concentration of the compound). For specific-important phenotypes, we also considered a  
731 specific phenotype to be conserved if a potential ortholog had fitness  $< -1$  and  $t < -4$  in a similar  
732 experiment.

733  
734 **“Predicting” TIGR subroles from cofitness or conserved cofitness.** We only considered  
735 hits that were in the top 10 cofit hits for a gene, only hits with cofitness above 0.4 (or conserved  
736 cofitness above 0.4), and only hits that were at least 10 kilobases from the query gene. In these  
737 cases, we predict that the gene has the same subrole as the best-scoring hit. When testing  
738 cofitness, the hits were sorted by the cofitness. When testing conserved cofitness, the hits were  
739 sorted by the lower of the cofitness in this organism and the best cofitness for orthologs in other  
740 organisms.

741  
742 **Genome and gene annotations.** For previously-published genomes, gene annotations were  
743 taken from MicrobesOnline, Integrated Microbial Genomes (IMG), or RefSeq. Newly-sequenced  
744 genomes were annotated with RAST<sup>49</sup>, except that *S. koreensis* DSMZ 15582 was annotated  
745 by IMG. See **Supplementary Table 16** for a summary of genome annotations and their  
746 accession numbers.

747  
748 **Classification of how informative annotations are.** To assess the existing computational  
749 annotations for these genomes, we classified all of their proteins into one of four groups:  
750 detailed TIGR role, hypothetical, vague, or other detailed. (1) “Detailed TIGR role” includes  
751 proteins that belong to a TIGRFAM role other than “Unclassified”, “Unknown function”, or  
752 “Hypothetical proteins” and had a subrole other than “Unknown substrate”, “Two-component  
753 systems”, “Role category not yet assigned”, “Other”, “General”, “Enzymes of unknown  
754 specificity”, “Domain”, or “Conserved”. (2) “Hypothetical” includes proteins containing the  
755 annotation description “hypothetical protein”, “unknown function”, “uncharacterized”, or if the  
756 entire description matched “TIGRnnnnn family protein” or “membrane protein”. (3) Proteins were  
757 considered to have “vague” annotations if the gene description contained “family”, “domain  
758 protein”, “related protein”, “transporter related”, or if the entire description matched common  
759 non-specific annotations (“abc transporter atp-binding protein”, “abc transporter permease”,  
760 “abc transporter substrate-binding protein”, “abc transporter”, “acetyltransferase”, “alpha/beta  
761 hydrolase”, “aminohydrolase”, “aminotransferase”, “atpase”, “dehydrogenase”, “dna-binding  
762 protein”, “fad-dependent oxidoreductase”, “gcn5-related n-acetyltransferase”, “histidine kinase”,  
763 “hydrolase”, “lipoprotein”, “membrane protein”, “methyltransferase”, “mfs transporter”,  
764 “oxidoreductase”, “permease”, “porin”, “predicted dna-binding transcriptional regulator”,  
765 “predicted membrane protein”, “probable transmembrane protein”, “putative membrane protein”,  
766 “response regulator receiver protein”, “rnd transporter”, “sam-dependent methyltransferase”,  
767 “sensor histidine kinase”, “serine/threonine protein kinase”, “signal peptide protein”, “signal  
768 transduction histidine kinase”, “tonb-dependent receptor”, “transcriptional regulator”,  
769 “transcriptional regulators”, or “transporter”). The remaining proteins were considered to have  
770 “other detailed” annotations.

771  
772 To identify a subset of the proteins annotated as “hypothetical” or “vague” that do not belong to  
773 any characterized families, we relied on Pfam and TIGRFAM. A Pfam was considered to be  
774 uncharacterized if its name began with either DUF or UPF (which is short for “uncharacterized  
775 protein family”). A TIGRFAM was considered to be uncharacterized if it had no link to a role or if  
776 the top-level role was “Unknown function”. To identify poorly-annotated proteins from diverse  
777 bacteria (for **Fig. 5a**), we used the rules for vague annotations only.

778

779 **Availability of data and code.** See the Fitness Browser (<http://fit.genomics.lbl.gov>) or the data  
780 downloads page (<http://morgannprice.org/bigfit/>), which includes the supplementary information.  
781 The BarSeq or TnSeq reads were analyzed with the RB-TnSeq scripts  
782 (<https://bitbucket.org/berkeleylab/feba>); we used statistics versions 1.0.2, 1.0.3, or 1.1.0 of the  
783 code.

## 784 Acknowledgements

785 We thank Axel Visel for help editing the manuscript and Victoria Lo, Wenjun Shao, and Keith  
786 Keller for technical assistance with the Fitness Browser web site. Sequencing was performed at:  
787 the Vincent J. Coates Genomics Sequencing Laboratory (University of California at Berkeley),  
788 supported by NIH S10 Instrumentation Grants S10RR029668, S10RR027303, and OD018174;  
789 at the DOE Joint Genome Institute; at the College of Biological Sciences <sup>UC</sup>DNA Sequencing  
790 Facility (UC Davis); and at the Institute for Genomics Sciences (University of Maryland).

791

792 Studies of novel isolates were conducted by ENIGMA and were supported by the Office of  
793 Science, Office of Biological and Environmental Research of the U.S. Department of Energy,  
794 under contract DE-AC02-05CH11231. The other data collection was supported by Laboratory  
795 Directed Research and Development (LDRD) funding from Berkeley Lab, provided by the  
796 Director, Office of Science, of the U.S. Department of Energy under contract DE-AC02-  
797 05CH11231 and a Community Science Project from the Joint Genome Institute to M.J.B., J.B.,  
798 A.P.A., and A.D. The work conducted by the U.S. Department of Energy Joint Genome Institute,  
799 a DOE Office of Science User Facility, is supported by the Office of Science of the U.S.  
800 Department of Energy under contract no. DE-AC02-05CH11231.

## 801 Author contributions

802 AMD, APA, MNP, MJB, and JB conceived the project. AMD, APA, MJB, and JB supervised the  
803 project. AMD led the experimental work. AMD, KMW, RJW, RAM, MC, JR, JVK, JSL, YS, ZE,  
804 and HS collected data. RC isolated bacteria. MNP and AMD analyzed the fitness data. RAM,  
805 RJW, and MNP assembled genomes. BER provided resources and advice on *S. elongatus*  
806 experiments. MNP, MJB, and AMD wrote the paper.

## 807 References

808

- 809 1. Schnoes, A. M., Brown, S. D., Dodevski, I. & Babbitt, P. C. Annotation error in public  
810 databases: misannotation of molecular function in enzyme superfamilies. *PLoS Comput.*  
811 *Biol.* **5**, e1000605 (2009).
- 812 2. Chang, Y.-C. *et al.* COMBEX-DB: an experiment centered database of protein function:  
813 knowledge, predictions and knowledge gaps. *Nucleic Acids Res.* **44**, D330–5 (2016).
- 814 3. Deutschbauer, A. *et al.* Towards an informative mutant phenotype for every bacterial  
815 gene. *J. Bacteriol.* **196**, 3643–3655 (2014).
- 816 4. Deutschbauer, A. *et al.* Evidence-based annotation of gene function in *Shewanella*  
817 *oneidensis* MR-1 using genome-wide fitness profiling across 121 conditions. *PLoS Genet.*  
818 **7**, e1002385 (2011).
- 819 5. Nichols, R. J. *et al.* Phenotypic landscape of a bacterial cell. *Cell* **144**, 143–156 (2011).
- 820 6. Price, M. N. *et al.* The genetic basis of energy conservation in the sulfate-reducing  
821 bacterium *Desulfovibrio alaskensis* G20. *Front Microbiol* **5**, 577 (2014).
- 822 7. Langridge, G. C. *et al.* Simultaneous assay of every *Salmonella* Typhi gene using one  
823 million transposon mutants. *Genome Res* **19**, 2308–2316 (2009).
- 824 8. van Opijnen, T., Bodi, K. L. & Camilli, A. Tn-seq: high-throughput parallel sequencing for  
825 fitness and genetic interaction studies in microorganisms. *Nat Methods* **6**, 767–772  
826 (2009).

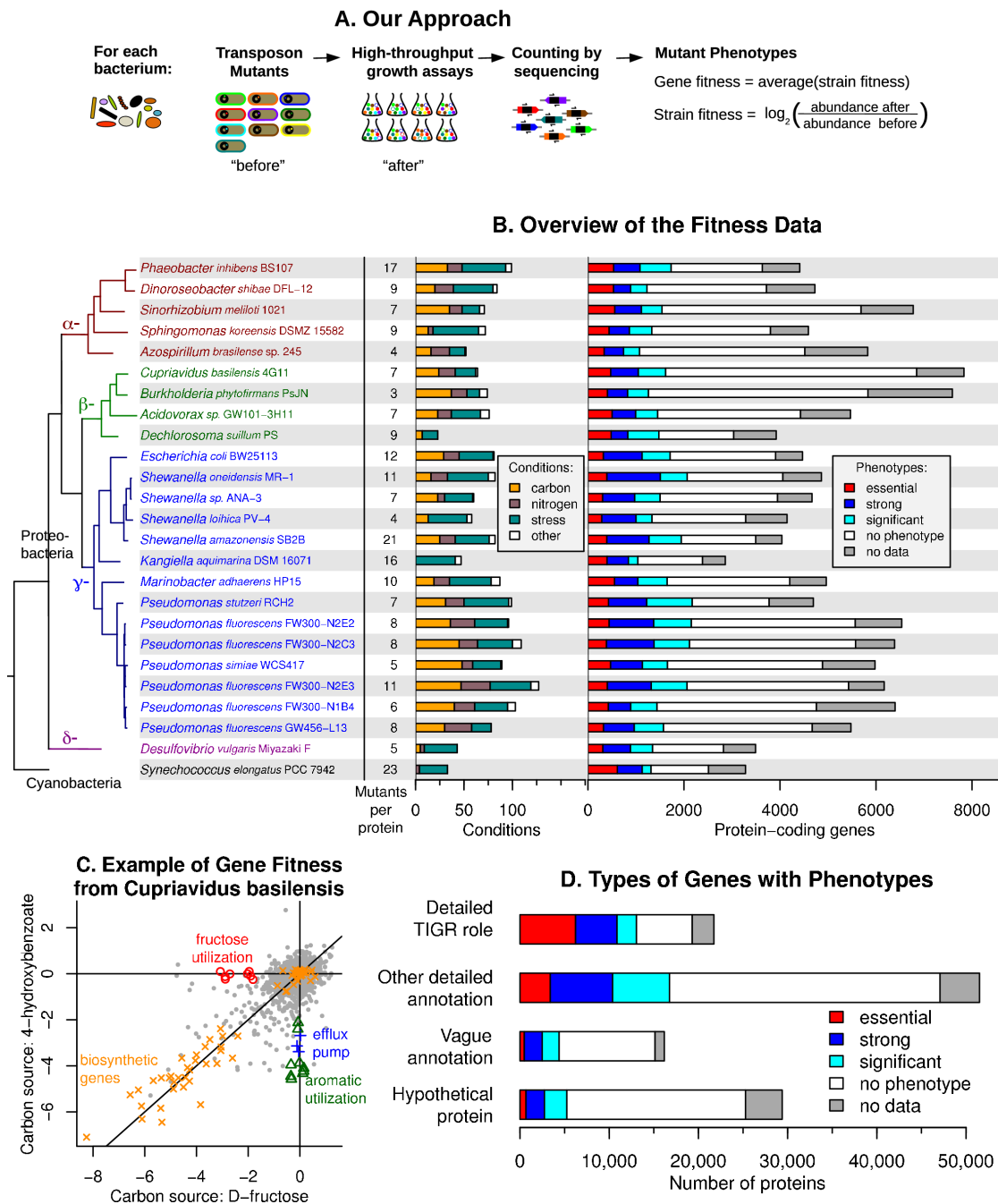


- 827 9. Wetmore, K. M. *et al.* Rapid quantification of mutant fitness in diverse bacteria by  
828 sequencing randomly bar-coded transposons. *MBio* **6**, e00306–15 (2015).
- 829 10. Rubin, B. E. *et al.* The essential gene set of a photosynthetic organism. *Proc. Natl. Acad.*  
830 *Sci. U.S.A.* **112**, E6634–43 (2015).
- 831 11. Melnyk, R. A. *et al.* Novel mechanism for scavenging of hypochlorite involving a  
832 periplasmic methionine-rich Peptide and methionine sulfoxide reductase. *MBio* **6**,  
833 e00233–15 (2015).
- 834 12. Smith, A. M. *et al.* Quantitative phenotyping via deep barcode sequencing. *Genome Res*  
835 **19**, 1836–1842 (2009).
- 836 13. Rensing, C., Pribyl, T. & Nies, D. H. New functions for the three subunits of the CzcCBA  
837 cation-proton antiporter. *J. Bacteriol.* **179**, 6871–6879 (1997).
- 838 14. Hottes, A. K. *et al.* Bacterial Adaptation through Loss of Function. *PLoS Genet.* **9**,  
839 e1003617 (2013).
- 840 15. Finn, R. D. *et al.* Pfam: the protein families database. *Nucleic Acids Res.* **42**, D222–30  
841 (2014).
- 842 16. Haft, D. H. *et al.* TIGRFAMs and Genome Properties in 2013. *Nucleic Acids Res.* **41**,  
843 D387–95 (2013).
- 844 17. Baker, J. L. *et al.* Widespread genetic switches and toxicity resistance proteins for  
845 fluoride. *Science* **335**, 233–235 (2012).
- 846 18. Keseler, I. M. *et al.* EcoCyc: fusing model organism databases with systems biology.  
847 *Nucleic Acids Res.* **41**, D605–12 (2013).
- 848 19. Rabus, R., Reizer, J., Paulsen, I. & Saier, M. H. Enzyme I(Ntr) from Escherichia coli. A  
849 novel enzyme of the phosphoenolpyruvate-dependent phosphotransferase system  
850 exhibiting strict specificity for its phosphoryl acceptor, NPr. *J Biol Chem* **274**, 26185–  
851 26191 (1999).
- 852 20. Justice, S. S., Hunstad, D. A., Cegelski, L. & Hultgren, S. J. Morphological plasticity as a  
853 bacterial survival strategy. *Nature Publishing Group* **6**, 162–168 (2008).
- 854 21. da Rocha, R. P., Paquola, A. C. de M., Marques, M. D. V., Menck, C. F. M. & Galhardo,  
855 R. S. Characterization of the SOS regulon of *Caulobacter crescentus*. *J. Bacteriol.* **190**,  
856 1209–1218 (2008).
- 857 22. Abella, M., Campoy, S., Erill, I., Rojo, F. & Barbé, J. Cohabitation of two different *lexA*  
858 regulons in *Pseudomonas putida*. *J. Bacteriol.* **189**, 8855–8862 (2007).
- 859 23. Cirz, R. T., O'Neill, B. M., Hammond, J. A., Head, S. R. & Romesberg, F. E. Defining the  
860 *Pseudomonas aeruginosa* SOS response and its role in the global response to the  
861 antibiotic ciprofloxacin. *J. Bacteriol.* **188**, 7101–7110 (2006).
- 862 24. Byrne, R. T., Chen, S. H., Wood, E. A., Cabot, E. L. & Cox, M. M. *Escherichia coli* genes  
863 and pathways involved in surviving extreme exposure to ionizing radiation. *J. Bacteriol.*  
864 **196**, 3534–3545 (2014).
- 865 25. Wiegmann, K. *et al.* Carbohydrate catabolism in *Phaeobacter inhibens* DSM 17395, a  
866 member of the marine roseobacter clade. *Appl Environ Microbiol* **80**, 4725–4737 (2014).
- 867 26. Brouns, S. J. J. *et al.* Identification of the missing links in prokaryotic pentose oxidation  
868 pathways: evidence for enzyme recruitment. *J Biol Chem* **281**, 27378–27388 (2006).
- 869 27. Stephens, C. *et al.* Genetic analysis of a novel pathway for D-xylose metabolism in  
870 *Caulobacter crescentus*. *J. Bacteriol.* **189**, 2181–2185 (2007).
- 871 28. Johnsen, U. *et al.* D-xylose degradation pathway in the halophilic archaeon *Haloferax*  
872 *volcanii*. *J Biol Chem* **284**, 27290–27303 (2009).
- 873 29. Tagourt, J., Landoulsi, A. & Richarme, G. Cloning, expression, purification and  
874 characterization of the stress kinase YeaG from *Escherichia coli*. *Protein Expr. Purif.* **59**,  
875 79–85 (2008).
- 876 30. Clark, W. T. & Radivojac, P. Analysis of protein function and its prediction from amino  
877 acid sequence. *Proteins* **79**, 2086–2096 (2011).
- 878 31. Markowitz, V. M. *et al.* IMG 4 version of the integrated microbial genomes comparative  
879 analysis system. *Nucleic Acids Res.* **42**, D560–7 (2014).
- 880 32. Dehal, P. S. *et al.* MicrobesOnline: an integrated portal for comparative and functional

- 881 genomics. *Nucleic Acids Res.* **38**, D396–400 (2010).
- 882 33. Thorgersen, M. P. *et al.* Molybdenum Availability is Key to Nitrate Removal in  
883 Contaminated Groundwater Environments. *Appl Environ Microbiol* AEM.00917–15  
884 (2015). doi:10.1128/AEM.00917-15
- 885 34. Ray, J. *et al.* Complete Genome Sequence of *Cupriavidus basilensis* 4G11, Isolated from  
886 the Oak Ridge Field Research Center Site. *Genome Announc* **3**, e00322–15 (2015).
- 887 35. Baba, T. *et al.* Construction of *Escherichia coli* K-12 in-frame, single-gene knockout  
888 mutants: the Keio collection. *Mol. Syst. Biol.* **2**, 2006 0008 (2006).
- 889 36. Kuehl, J. V. *et al.* Functional genomics with a comprehensive library of transposon  
890 mutants for the sulfate-reducing bacterium *Desulfovibrio alaskensis* G20. *MBio* **5**,  
891 e01041–14 (2014).
- 892 37. Zane, G. M., Yen, H. C. & Wall, J. D. Effect of the deletion of *qmoABC* and the promoter-  
893 distal gene encoding a hypothetical protein on sulfate reduction in *Desulfovibrio vulgaris*  
894 Hildenborough. *Appl Environ Microbiol* **76**, 5500–5509 (2010).
- 895 38. Kahm, M., Hasenbrink, G., Lichtenberg-Frate, H., Ludwig, J. & Kschischo, M. grofit:  
896 Fitting Biological Growth Curves with R. *Journal of Statistical Software* **33**, (2010).
- 897 39. Bankevich, A. *et al.* SPAdes: a new genome assembly algorithm and its applications to  
898 single-cell sequencing. *J. Comput. Biol.* **19**, 455–477 (2012).
- 899 40. Bashir, A. *et al.* A hybrid approach for the automated finishing of bacterial genomes. *Nat.*  
900 *Biotechnol.* **30**, 701–707 (2012).
- 901 41. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat Methods*  
902 **9**, 357–359 (2012).
- 903 42. Walker, B. J. *et al.* Pilon: an integrated tool for comprehensive microbial variant detection  
904 and genome assembly improvement. *PLoS One* **9**, e112963 (2014).
- 905 43. Tritt, A., Eisen, J. A., Facciotti, M. T. & Darling, A. E. An integrated pipeline for de novo  
906 assembly of microbial genomes. *PLoS One* **7**, e42304 (2012).
- 907 44. Hunt, M. *et al.* Circlator: automated circularization of genome assemblies using long  
908 sequencing reads. *Genome Biol* **16**, 294 (2015).
- 909 45. Eddy, S. R. Accelerated Profile HMM Searches. *PLoS Comput. Biol.* **7**, e1002195 (2011).
- 910 46. Lerat, E., Daubin, V., Ochman, H. & Moran, N. A. Evolutionary origins of genomic  
911 repertoires in bacteria. *PLoS Biol* **3**, e130 (2005).
- 912 47. Wu, M. & Scott, A. J. Phylogenomic analysis of bacterial and archaeal sequences with  
913 AMPHORA2. *Bioinformatics* **28**, 1033–1034 (2012).
- 914 48. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2--approximately maximum-likelihood  
915 trees for large alignments. *PLoS One* **5**, e9490 (2010).
- 916 49. Aziz, R. K. *et al.* The RAST Server: rapid annotations using subsystems technology. *BMC*  
917 *Genomics* **9**, 75 (2008).
- 918
- 919
- 920
- 921
- 922
- 923
- 924

925 **Figures and legends**

926



927

928

929

930

931

932

933

934

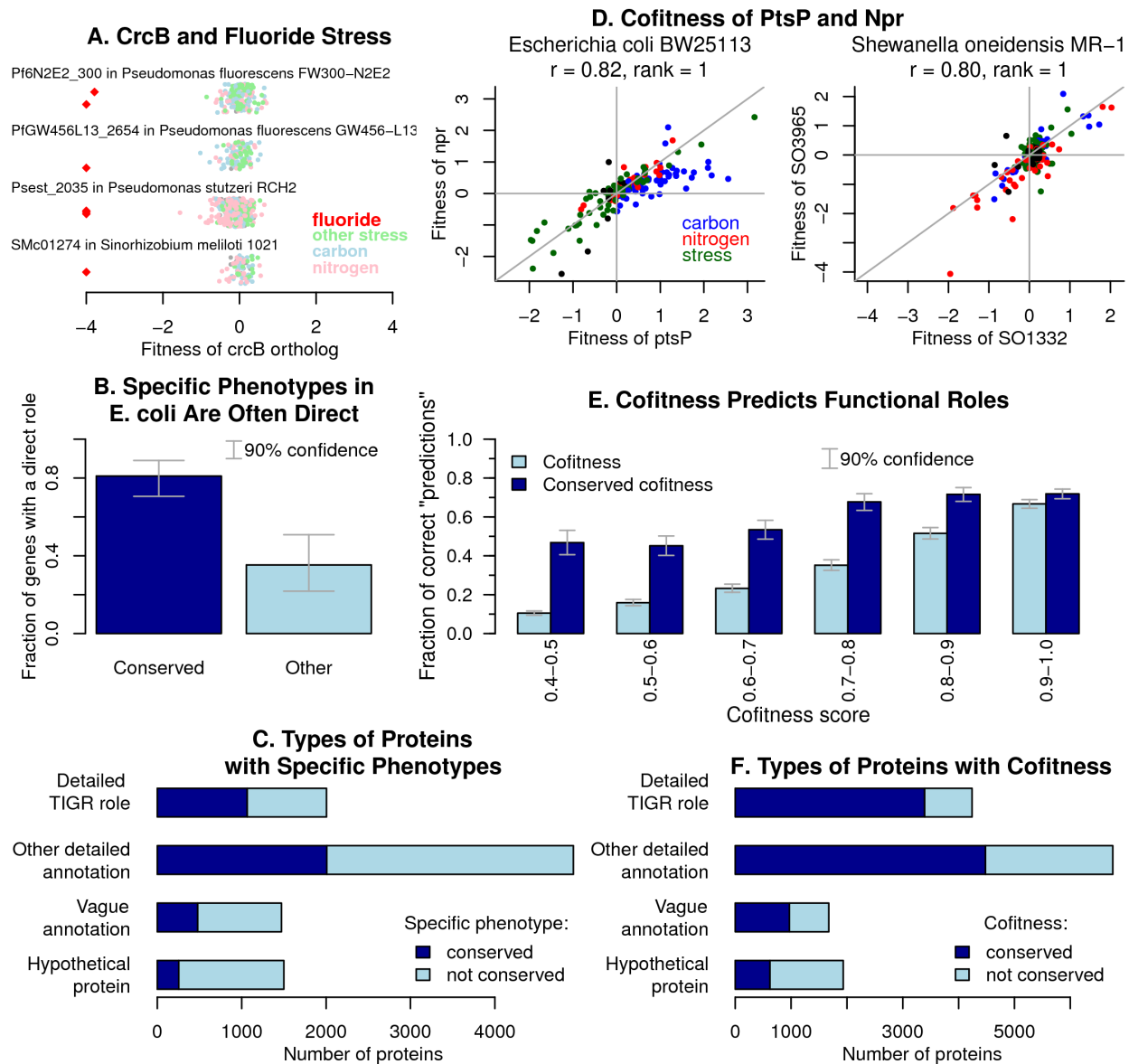
935

936

937

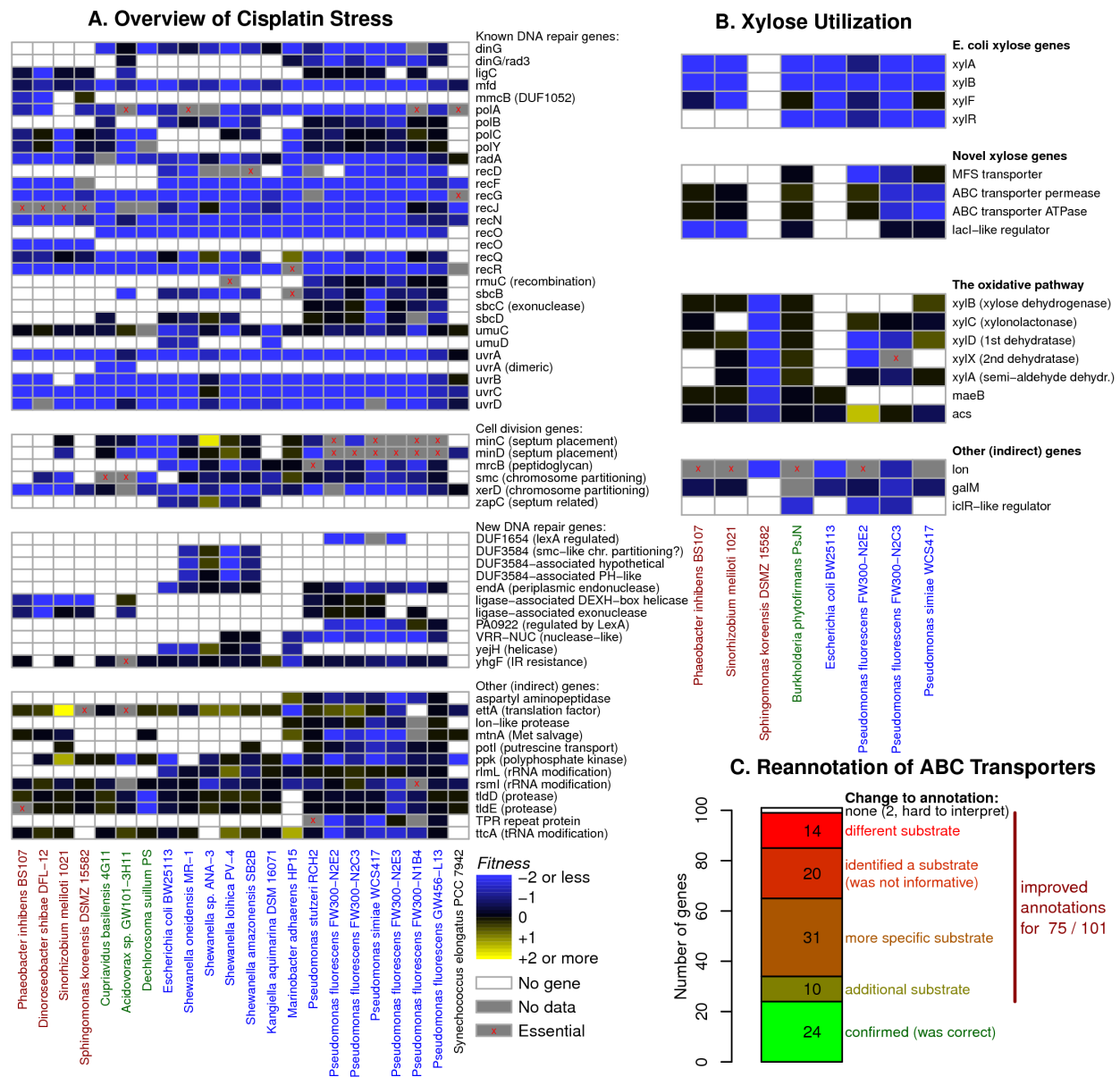
938

**Figure 1. High-throughput genetics for 25 bacteria.** (A) Our approach for measuring gene fitness. (B) Overview of our data. For each bacterium, we show the number of mutant strains that we used to estimate fitness for a typical protein (Methods), the types of conditions that we studied, and how many proteins had mutant phenotypes. A strong phenotype is defined as fitness < -2. (C) Gene fitness during the utilization of two carbon sources by *Cupriavidus basilensis*. See **Supplementary Table 5** for details on the highlighted genes. The 4-hydroxybenzoate data is the average from two biological replicates. (D) We classified the proteins from all 25 bacteria by how informative their annotations were (Methods), and for each class, we show how many proteins have each type of phenotype.

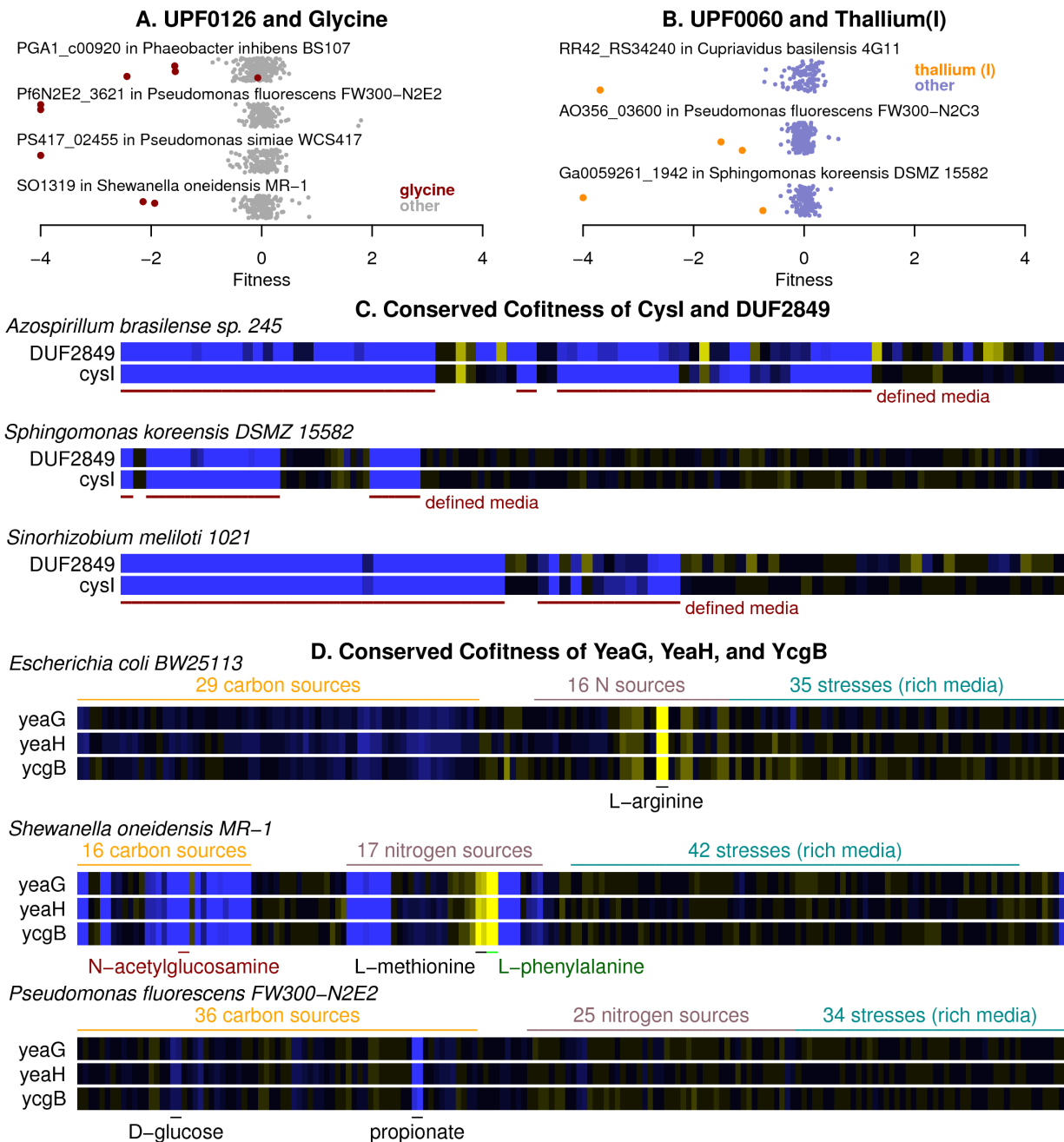


939  
 940 **Figure 2. Identification of conserved phenotypes.** (A) An example of a conserved and  
 941 specific phenotype. Each point shows the fitness of *crcB* in an experiment, with fluoride stress  
 942 experiments highlighted in red. Values less than -4 are shown at -4. The y-axis is random. (B)  
 943 The fraction of *E. coli* proteins with a specific phenotype in defined media that are 'directly'  
 944 involved in the uptake or catabolism of the compound. Conserved indicates that an ortholog of  
 945 the *E. coli* protein from another bacterium is important for fitness during growth with the same  
 946 compound (fitness < -1). (C) How many proteins of each type had a conserved specific  
 947 phenotype or a specific phenotype. (D) Comparison of fitness values for *ptsP* and *npr* from  
 948 *E. coli* and *S. oneidensis* across all tested conditions. The experiments are color coded by type.  
 949 (E) Using TIGR subroles to test the accuracy of the gene-gene associations. (F) How many  
 950 proteins of each type had at least one association from conserved cofilness ( $r > 0.6$  in both  
 951 bacteria) or else from cofilness ( $r > 0.8$ ).  
 952



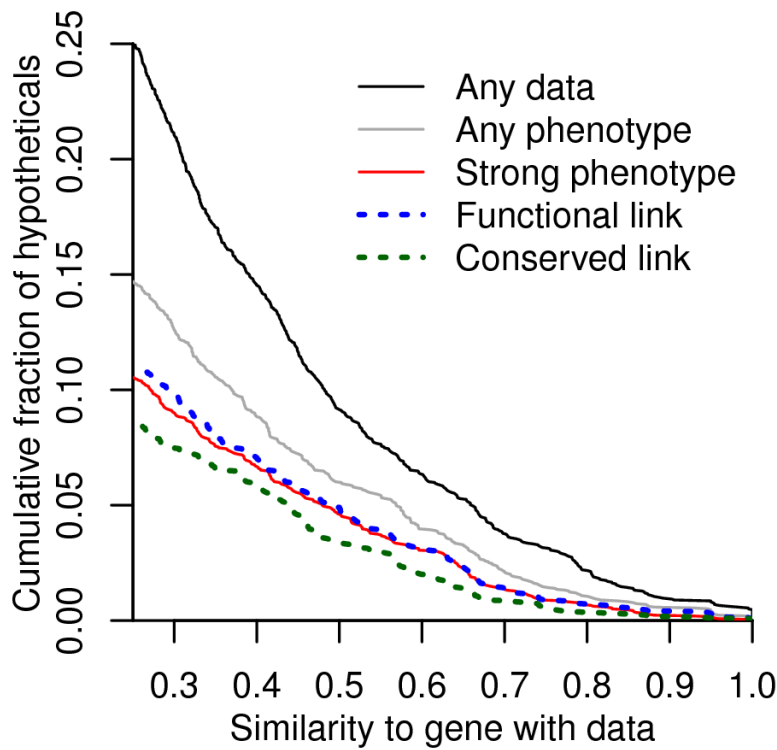


953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966



967  
 968  
 969  
 970  
 971  
 972  
 973  
 974  
 975

**Figure 4. Annotation of proteins from uncharacterized families.** (A, B) Conserved specific phenotypes for proteins of unknown function. Each point represents the fitness of the protein in an individual experiment. Values under -4 are shown at -4. The y-axis is random. (C) Heatmap of fitness data for *cysI* (sulfite reductase) and DUF2849 in three bacteria. (D) Heatmap of fitness data for *yeaG*, *yeaH*, and *ycgB* from three different bacteria. Certain experimental conditions with significant phenotypes are highlighted. The color scale is the same as in Fig. 3a.



976  
977

978 **Figure 5. Relevance to all bacteria.** We selected hypothetical or vaguely-annotated proteins  
979 from diverse bacterial species, we compared them to the genes that we have fitness data for  
980 (using protein BLAST), and we identified potential orthologs as best hits that were homologous  
981 over at least 75% of each protein's length. We show the fraction of these proteins that have an  
982 ortholog with each type of phenotype and that is above a given level of sequence similarity.  
983 Similarity was defined as the ratio of the alignment's bit score to the score from aligning the  
984 query to itself.  
985

## Supplementary information for “Deep Annotation of Protein Function across Diverse Bacteria from Mutant Phenotypes”

Morgan N. Price<sup>1</sup>, Kelly M. Wetmore<sup>1</sup>, R. Jordan Waters<sup>2</sup>, Mark Callaghan<sup>1</sup>, Jayashree Ray<sup>1</sup>, Jennifer V. Kuehl<sup>1</sup>, Ryan A. Melnyk<sup>1</sup>, Jacob S. Lamson<sup>1</sup>, Yumi Suh<sup>1</sup>, Zuelma Esquivel<sup>1</sup>, Harini Sadeeshkumar<sup>1</sup>, Romy Chakraborty<sup>3</sup>, Benjamin E. Rubin<sup>4</sup>, James Bristow<sup>2</sup>, Matthew J. Blow<sup>2,\*</sup>, Adam P. Arkin<sup>1,5,\*</sup>, Adam M. Deutschbauer<sup>1,\*</sup>

<sup>1</sup>Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory

<sup>2</sup>Joint Genome Institute, Lawrence Berkeley National Laboratory

<sup>3</sup>Climate and Ecosystem Sciences Division, Lawrence Berkeley National Laboratory

<sup>4</sup>Division of Biological Sciences, University of California, San Diego

<sup>5</sup>Department of Bioengineering, University of California, Berkeley

\*To whom correspondence should be addressed:

MJB ([MJBlow@lbl.gov](mailto:MJBlow@lbl.gov))

APA ([APArkin@lbl.gov](mailto:APArkin@lbl.gov))

AMD ([AMDeutschbauer@lbl.gov](mailto:AMDeutschbauer@lbl.gov))

Website for interactive analysis of mutant fitness data:

<http://fit.genomics.lbl.gov/>

Website with supplementary information and bulk data downloads:

<http://genomics.lbl.gov/supplemental/bigfit/>

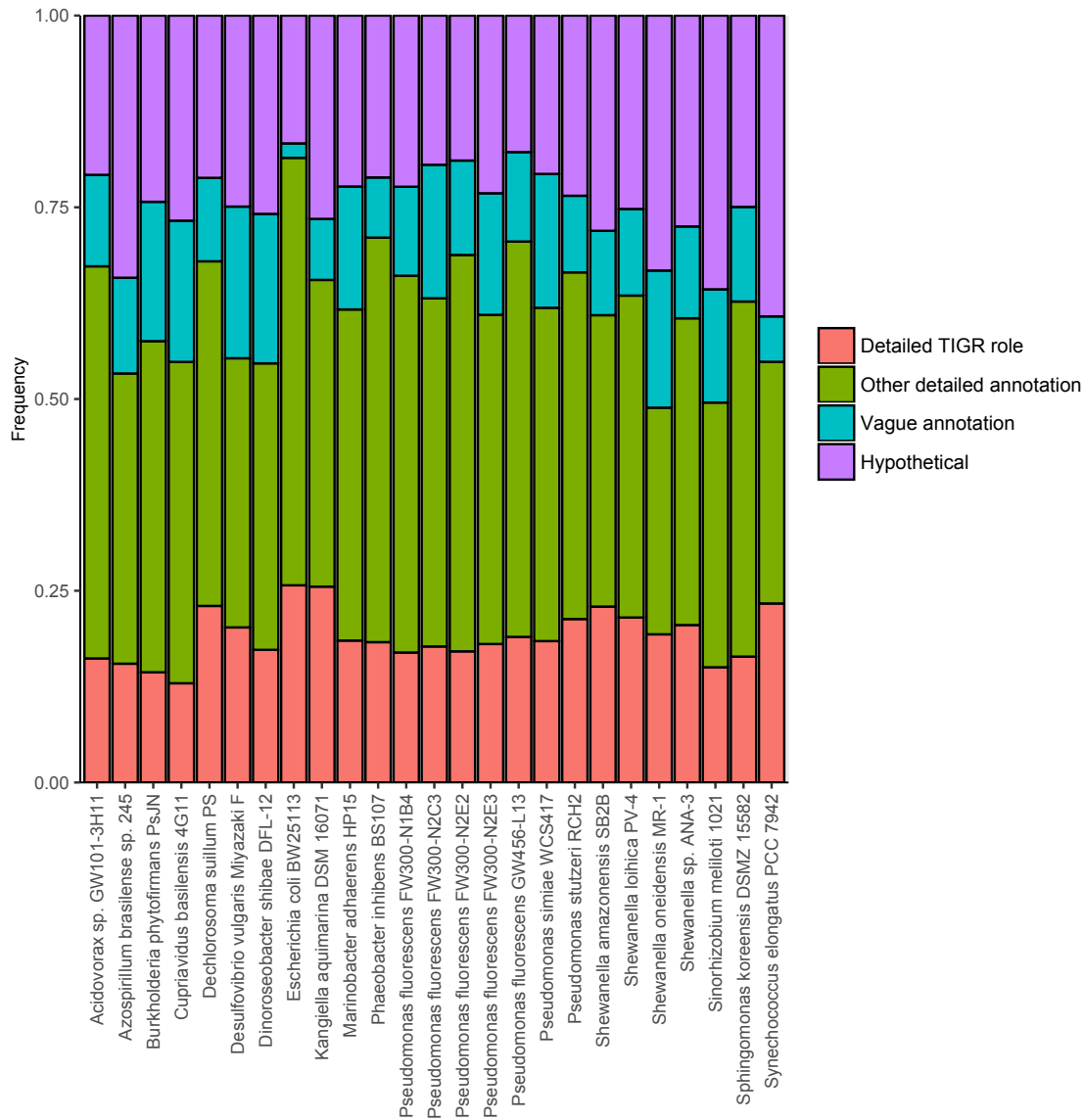


## Contents

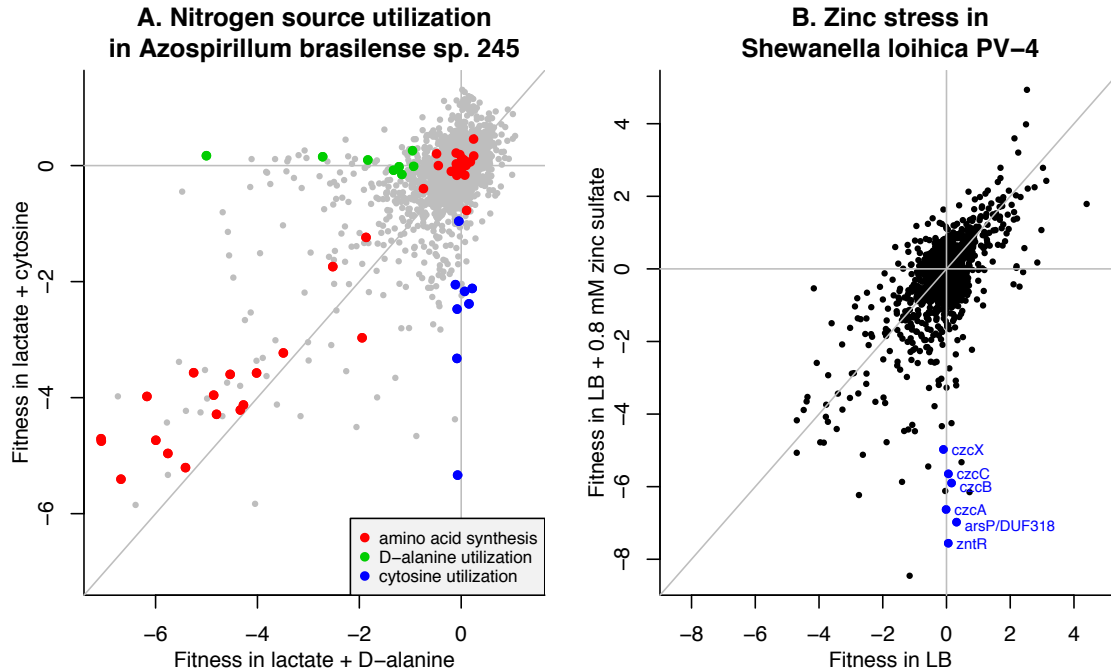
<b>Page 3</b>	Supplementary Figure 1 – Protein annotations for the 25 bacteria
<b>Page 4</b>	Supplementary Figure 2 – Examples of nitrogen source and stress fitness experiments
<b>Page 5</b>	Supplementary Figure 3 – Phenotypes versus types of genes
<b>Page 6</b>	Supplementary Figure 4 – EndA is important for cisplatin resistance
<b>Page 7</b>	Supplementary Figure 5 – Growth of signaling mutants in <i>Escherichia coli</i>
<b>Page 8</b>	Supplementary Figure 6 – Growth of signaling mutants from <i>Shewanella oneidensis</i>
<b>Page 9</b>	Supplementary Note 1 – Validation of essential genes analysis
<b>Page 11</b>	Supplementary Note 2 – Rationales for annotating domains of unknown function

### In separate excel file:

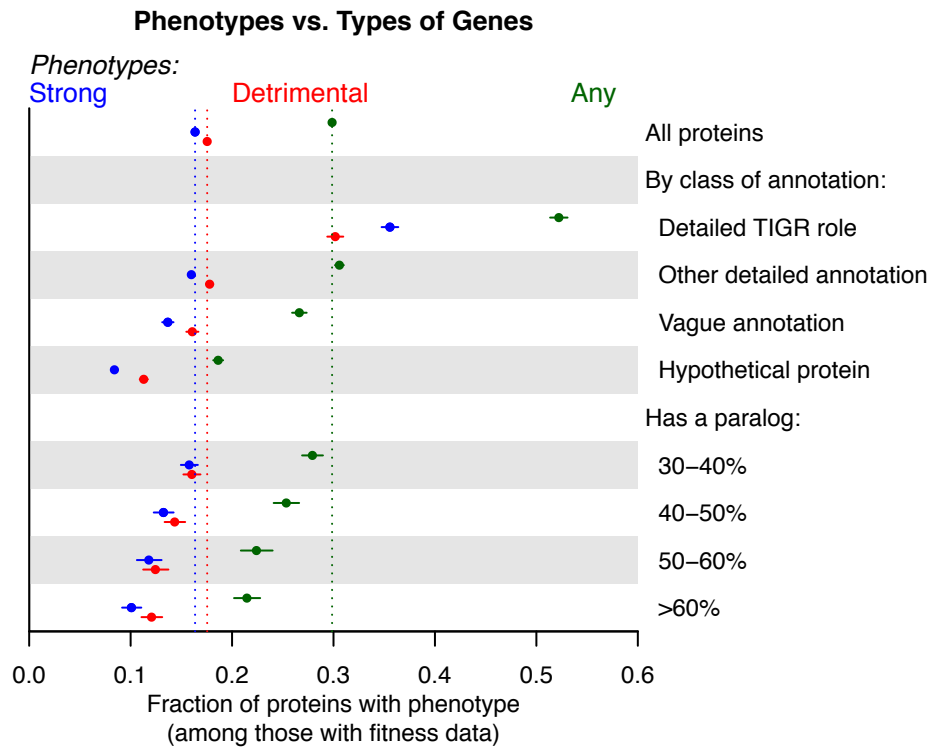
Supplementary Table 1 – List of essential genes
Supplementary Table 2 – Growth on carbon substrates
Supplementary Table 3 – Growth on nitrogen substrates
Supplementary Table 4 – Growth on stress compounds
Supplementary Table 5 – Catabolic genes in <i>Cupriavidus basilensis</i>
Supplementary Table 6 – <i>E. coli</i> genes with specific phenotypes in carbon and nitrogen source experiments
Supplementary Table 7 – Genes with conserved specific phenotypes or conserved cofitness
Supplementary Table 8 – Genes with conserved specific phenotypes under cisplatin stress
Supplementary Table 9 – Genes with specific phenotypes during D-xylose utilization
Supplementary Table 10 – Reannotation of ABC transporters
Supplementary Table 11 – Uncharacterized protein families with conserved specific phenotypes or conserved cofitness
Supplementary Table 12 – Bacteria used in this study
Supplementary Table 13 – Media formulations used in this study
Supplementary Table 14 – Genome sequencing statistics
Supplementary Table 15 – Transposon mutagenesis details
Supplementary Table 16 – Genome annotation summary



**Supplementary Figure 1: Protein annotations for the 25 bacteria.** The fraction of proteins in each genome with different levels of annotation (see main text and Methods) is color-coded.

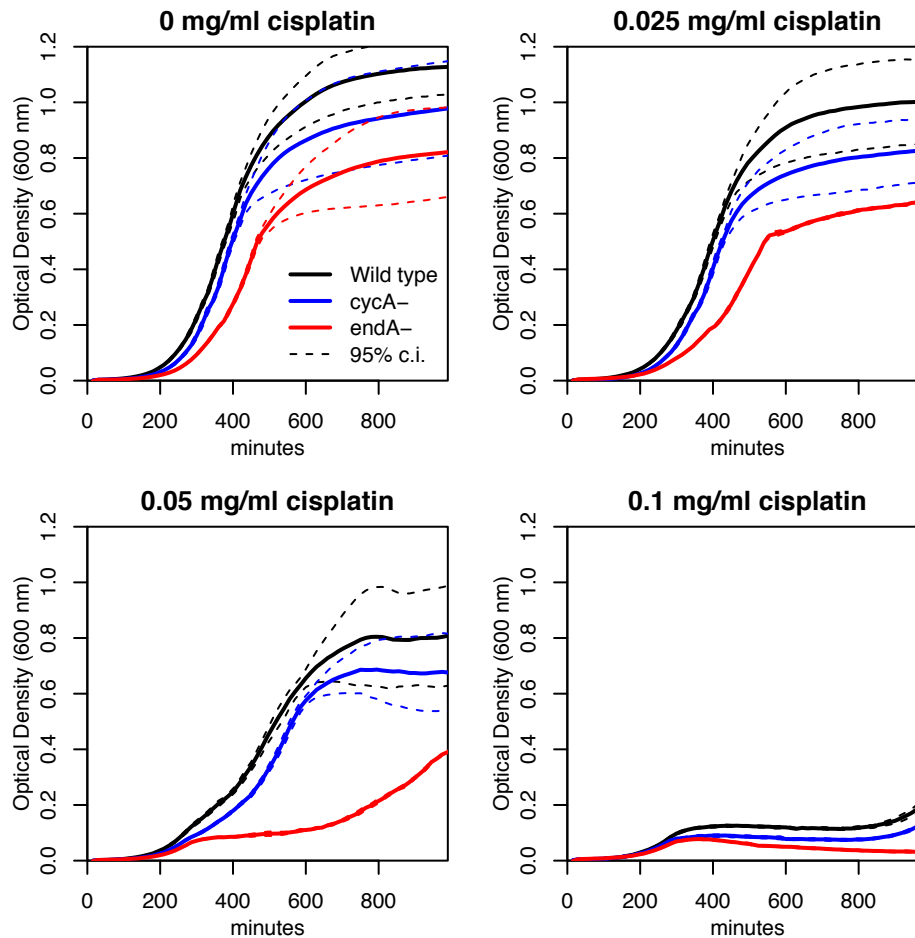


**Supplementary Figure 2: Examples of nitrogen source and stress fitness experiments.** (A) The utilization of D-alanine or cytosine by *Azospirillum brasilense* sp. 245. Each point shows the fitness of a gene in the two conditions. The data is the average of two biological replicates for each nitrogen source. Amino acid synthesis genes were identified using the top-level role in TIGRFAMs. The genes for D-alanine utilization were a D-amino acid dehydrogenase (AZOBR\_RS08020), an ABC transporter operon (AZOBR\_RS08235:RS08260), and a LysR family regulator (AZOBR\_RS21915). The genes for cytosine utilization were cytosine deaminase (AZOBR\_RS31895) and two ABC transporter operons (AZOBR\_RS06950:RS06965 and AZOBR\_RS31875:RS31885). (B) Zinc stress in *Shewanella loihica* PV-4. We compare fitness in rich media with added zinc (II) sulfate to fitness in plain rich media. The LB data is the average of two biological replicates. The highlighted genes include a putative heavy metal efflux pump (CzcCBA or Shew\_3358:Shew\_3356), a hypothetical protein at the beginning of the *czc* operon (CzcX), a zinc-responsive regulator (ZntR or Shew\_3411), and another heavy metal efflux gene related to *arsP* or DUF318 (Shew\_3410). *CzcX* lacks homology to any characterized protein, but homologs in other strains of *Shewanella* are also specifically important for resisting zinc stress. In both panels, the lines show  $x = 0$ ,  $y = 0$ , and  $x = y$ .

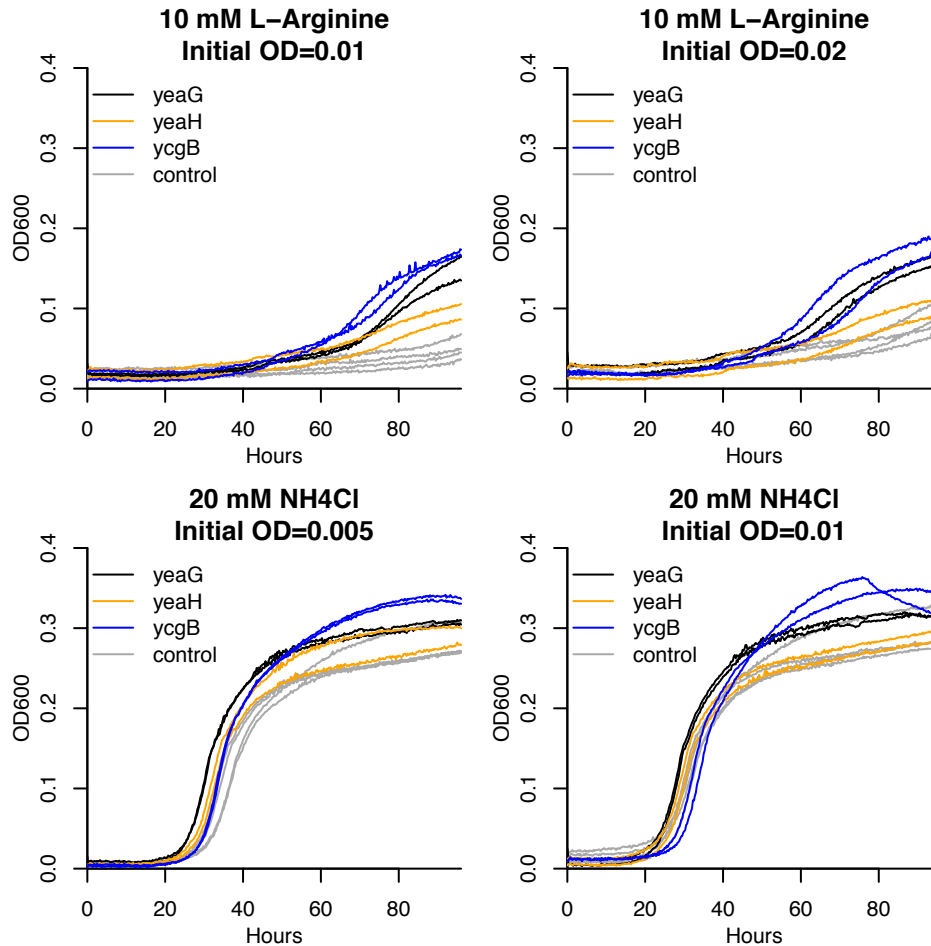


**Supplementary Figure 3. Phenotypes versus types of genes.** We categorized proteins in our data set by their type of annotation or by whether they have homologs in the same genome (“paralogs”). For each category, we show the fraction of genes that have significant phenotypes, and more specifically the fractions that have strong phenotypes (fitness < -2 and  $t < -5$ ) or are detrimental to fitness (fitness > 0).

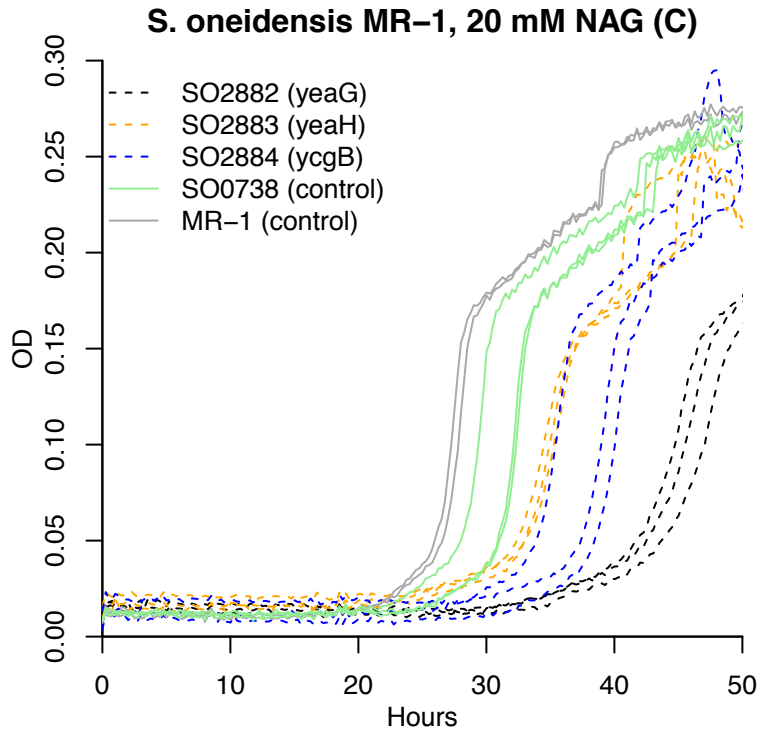




**Supplementary Figure 4. EndA is important for cisplatin resistance.** We compared the growth of the *endA*- deletion strain from the Keio collection to the growth of the base strain (BW25113) and the growth of a *cycA*- deletion strain (*cycA* encodes an amino acid transporter and is not expected to have a phenotype in this condition.) All experiments were conducted at 30°C in LB media with varying levels of cisplatin added. This experiment was conducted in a 96-well plate. Each growth curve is the average of 12 replicate wells and the dashed lines show 95% confidence intervals. *endA* encodes endonuclease I, which is believed to be in the periplasm -- it is released from cells after various shocks<sup>1,2</sup>. There are also mutants of *E. coli* that leak several periplasmic enzymes including some endonuclease I<sup>3</sup>. EndA has a likely signal peptide and no other transmembrane helix, which is consistent with being in the periplasm. It remains unclear how EndA plays a role in the response to DNA damage if it is located in the periplasm.



**Supplementary Figure 5. Growth of signaling mutants in *Escherichia coli*.** We grew deletion strains from the Keio collection<sup>4</sup> in M9 glucose media with varying nitrogen sources. The signaling mutants are in *yeaG*, *yeaH*, and *ycgB*. Control mutants are deletions of two pseudogenes, *agaA* or *ygaY*. The signaling mutants had a strong growth advantage on L-arginine, as expected from fitness assays with 10 mM L-arginine as the nitrogen source.



**Supplementary Figure 6. Growth of signaling mutants from *Shewanella oneidensis*.** We grew various strains derived from *S. oneidensis* MR-1 in a defined medium (ShewMM\_noCarbon) with 20 mM N-acetylglucosamine (NAG) as the carbon source. This medium contains ammonium as the nitrogen source. These mutants have transposons inserted within the signaling pathway (SO2882:SO2884) or in a pseudogene (SO0738, identified as such by Romine and colleagues<sup>5</sup>). We also grew wild-type *S. oneidensis* MR-1 as an additional control. Insertions in the signaling pathway had an increased lag and a slower growth rate, as expected from fitness assays in this condition. The mutants are from a previously described collection<sup>6</sup> and were generated independently from the mutants that were used to generate the fitness data.

## Supplementary Note 1. Validation of essential genes analysis

Our approach to identify essential genes was initially validated for *Synechococcus elongatus*<sup>7</sup>. To verify that this approach gave reasonable results for other bacteria, we compared our list of putatively essential proteins in *E. coli* K-12 to genes that are essential<sup>8</sup> or important for growth in LB<sup>4</sup>. We also compared the list of essential proteins in *Shewanella oneidensis* MR-1 to a previous analysis based on a different set of transposon mutants<sup>6</sup>. Finally, for all of the genomes, we examined the functional categories according to TIGRFAMs<sup>9</sup> and whether essentiality was conserved.

Of the 259 essential proteins in *E. coli*, 219 (85%) were in our list of essential proteins. Our list also included another 122 proteins (from 3,887 non-essential proteins), which corresponds to a false-discovery rate of 36%. Of those 122, 25 have reduced growth in LB, with OD<sub>600</sub> after 22 hours being in the bottom 10% of genes<sup>4</sup>. If we consider these as important genes then the FDR drops to 29%. The remaining false positives tend to be somewhat shorter than true essential genes, with median lengths of 638 and 942 nucleotides, respectively. These 97 “false positives” include 10 genes that are involved in protein synthesis or protein fate as well as other genes that are likely to be important for fitness (*minE*, *mviN*, *dapF*, *hoID*, *ubiB*, *ubiD*, and *purB*).

In *S. oneidensis* MR-1, our list contained 397 proteins. We compared this to a list of essential genes that we had previously generated using a different transposon and Sanger sequencing<sup>6</sup>. The previous analysis was conservative as genes that were not expected to be essential (based on orthology to *E. coli* or *Acinetobacter*, which are also  $\gamma$ -Proteobacteria) were required to be adjacent to another essential gene or to be conserved across most other *Shewanella* genomes. Of the 397 proteins in our new list, 298 were previously classified as essential and 83 were classified as unknown, Just 16 were previously identified as dispensable, and 2 of these are genes that are essential in *E. coli*. This implies a false discovery rate of around  $16/(397-83) = 5\%$ , with the caveat that insertions in a gene might be selected against even though the gene is not itself important (i.e., polar effects). We also examined the 15 proteins that are putatively essential in *S. oneidensis* MR-1 but lacked clear orthologs in the closely related strain *S. sp.* ANA-3, as these might be more likely to be false positives. Of these 15, at least seven are plausibly essential, including three prophage repressors that are probably required to prevent prophage excision; a gene adjacent to one of these repressors; the RepA protein that is probably required to maintain the megaplasmid, a tRNA synthetase (SO3128.2) with a putative internal stop codon that nevertheless forms full-length protein<sup>10</sup>, and ribosomal protein L25 (whose ortholog was missed because it seems to be annotated with the wrong start codon).

We then considered TIGR roles that are likely to be essential. The roles we chose were: DNA metabolism; transcription; protein synthesis; protein fate; energy metabolism; cell envelope; fatty acid and phospholipid metabolism; purines, pyrimidines, nucleosides,



and nucleotides; amino acid biosynthesis; and biosynthesis of cofactors, prosthetic groups, and carriers. We included biosynthetic and energetic genes as likely-essential roles because many bacteria cannot take up as wide a range of nutrients as *E. coli* can or have more limited ways of creating energy. In *E. coli* and *S. oneidensis*, these categories account for 62-67% of putatively essential genes but just 12-18% of other genes. In the 25 bacteria, these categories accounted for 44-69% of putatively essential genes but just 8-17% of other genes. This confirms that in each bacterium, most of these genes are essential.

Finally, we asked whether the putatively essential proteins had orthologs in other bacteria that were essential. Overall, 84% of the essential proteins were confirmed by conservation (they had an ortholog that was also essential). The organisms with the lowest proportions of conserved essentials were *S. elongatus* and *D. vulgaris* (56% and 64%, respectively). This might reflect their unique sources of energy (photosynthesis and dissimilatory sulfate reduction, respectively) or their evolutionary distance from the other bacteria that we studied.

## Supplementary Note 2. Rationales for annotating domains of unknown function

### Specific Predictions:

#### **DUF485 (PF04341): component of actP-like carboxylate transporters**

Several representatives of DUF485 (which is known as *yjcH* in *E. coli*) are cofit with an adjacent ActP-like permease. Examination of the per-strain data did not show any evidence of polar effects, so we suggest that DUF485 is required for the activity of the permease. The paper that characterized ActP in *E. coli* did not rule out the requirement of another gene <sup>11</sup>. DUF485 seems to be a membrane protein (i.e., AZOBR\_RS02935 has two transmembrane helices and a possible signal sequence). Together, this suggests that DUF485 is a component of the transporter. The clearest phenotypes are for pyruvate utilization (i.e. AZOBR\_RS02935).

#### **DUF1513 (PF07433): outer membrane component of ferrous iron uptake**

DUF1513 is in a conserved operon and is cofit with the other genes in that operon in multiple bacteria. The operon contains EfeO-like, DUF1111, a second EfeO-like, and DUF1513. EfeO has an unknown role in iron uptake by *efeUOB* and is a putative periplasmic lipoprotein. DUF1111 is proposed to be a homolog of EfeB, a di-heme peroxidase involved in ferrous iron uptake <sup>12</sup>. Thus, the operon seems to be involved in ferrous iron uptake. Related operons in  $\alpha$ -Proteobacteria often contain bacterioferritin, which is consistent with that role.

DUF1513 has a putative signal peptide (PSORTb) and has similarity to beta propellers with 6-8 repeats (Pfam clans). This suggests that it is the outer membrane component of this system.

Although the fitness data shows that DUF1513 is involved in this process, it is not certain that ferrous iron is the substrate. Several representatives of DUF1513 are pleiotropic, but AO356\_18450 is specifically important for chlorite resistance and Psest\_1156 is sensitive to various metals. Inhibiting iron uptake would plausibly create these phenotypes.

#### **DUF1656 (PF07869): component of efflux pump with MFP and FUSC**

Several representatives of DUF1656 are in an operon with and cofit with a RND efflux pump and a fusaric acid resistance-like protein. DUF1656 is related to *E. coli* YdhI and AaeX/YhcR. Homologs of these proteins are sometimes annotated as inner membrane efflux pump components or as Na<sup>+</sup>-dependent SNF-like transporters but we could not find the rationale for these annotations. (As of March 2016, EcoCyc reports that the functions of *aaeX* and *ydhI* are not known.)

This family has a variety of stress sensitivity phenotypes. In *Cupriavidus basilensis*, this efflux pump is specifically important for the utilization of 4-hydroxybenzoate. It could be involved in uptake, or 20 mM of 4-hydroxybenzoate may have been inhibitory and it is involved in efflux. Similarly, a number of representatives are important for the utilization of octanoate and it is not clear if this reflects uptake or efflux. In several strains of *Pseudomonas fluorescens*, the efflux pump is either important for or detrimental during acetate utilization. Finally, in *Zymomonas mobilis*, we found that a similar system (ZMO1432-ZMO1431-ZMO1430) is involved in resisting hydrolysate<sup>13</sup>.

#### **DUF1854 (PF08909): subunit of transporter for efflux of an amino acid polymer**

In two  $\beta$ -Proteobacteria, *Acidovorax* sp. 3H11 and *Cupriavidus basilensis* 4G11, representatives of DUF1854 (i.e., RR42\_RS04420) are cofit with two nearby genes that are annotated as cyanophycin synthetases as well as an ABC-like transporter. Cyanophycin is a copolymer of aspartate and arginine that many cyanobacteria use to store nitrogen; it is also formed by some heterotrophic bacteria<sup>14</sup>. The putative cyanophycin synthetases (ChpA and ChpA') have been studied in another strain of *Cupriavidus necator* and do not form cyanophycin but may produce a different light-scattering polymer<sup>15</sup>. Also, the genomes of many  $\beta$ -Proteobacteria contain *chpA* and *chpA'* but do not appear to encode the cyanophycinase (*chpB*) to break down the polymer, which hints at a different role for these genes<sup>14</sup>. As *chpA* and *chpA'* were cofit with an ABC transporter in both organisms, we propose that the polymer is being exported to form part of the cell wall rather than serving as a storage compound. This can also explain why the genes are important for motility (in *Acidovorax*) or have pleiotropic stress phenotypes (in both organisms). (We did not succeed in measuring fitness during motility in *C. basilensis*.) The ABC transporter contains an ABC transporter transmembrane domain and an ATPase domain. DUF1854 is usually found in an apparent operon with the ABC transporter, but in a few genomes the two proteins are fused together, as pointed out by the Pfam curators. Together with the cofitness, this suggests that DUF1854 is involved in the transport and forms a complex with the ABC transporter.

#### **DUF2849 (PF11011): electron source for sulfite reductase**

This family is usually upstream of *cysI*, the beta subunit of sulfite reductase. These sulfite reductases are important for fitness in our defined media, which confirms that they are indeed sulfite reductase and not nitrite reductase, as sulfate is the sulfur source in these media and sulfate must be reduced to sulfite and then to sulfide before it is assimilated. (In *S. meliloti*, *cysI* or SMc02124 was misannotated as "nitrite reductase.") However these genomes do not seem to contain the alpha flavoprotein subunit (*cysJ*). Instead, DUF2849 is found upstream. Several representatives of DUF2849 have similar fitness patterns as the downstream *cysI* (SMc01054, AZOBR\_RS10130, Ga0059261\_1497). The phenotypes of DUF2849 do not seem to be due to polar effects, as strains with insertions in either orientation have low fitness in defined media. Also,

DUF2849 is found fused to *cysI* in *Pseudomonas putida* GB-1. Altogether this suggests that DUF2849 is required for the activity of sulfite reductase. Usually *cysJ* is the electron source for *cysI* so we propose that in its absence, DUF2849 fulfills this role. (The relationship between the representatives of DUF2849 with similar cofitness was missed by the automated analysis of orthologs, as these alignments just missed the cutoffs of coverage > 80% or  $E < 10^{-5}$ . So this case is absent from **Supplementary table 11.**)

#### **DUF4212 (PF13937, TIGR03647): small subunit of transporter for D-alanine, lactate**

DUF4212 was predicted by the TIGRFAM curators to be the small subunit of a solute:sodium symporter because of its hydrophobicity and conserved gene context. Multiple members of this family (e.g., Sama\_1522 or Psest\_0346) are cofit with a putative large subunit of a symporter that is downstream (e.g., Sama\_1523 or Psest\_0347). Sama\_1522 from *Shewanella amazonensis* SB2B is important for fitness when D,L-lactate is the carbon source; similarly, a homolog in *S. oneidensis* MR-1 (SO2858) is more mildly important in some D,L-lactate conditions. Psest\_0346 in *Pseudomonas stutzeri* RCH2 (51% identical to Sama\_1523) is very important for D-alanine utilization and strongly detrimental during L-alanine utilization. Note that alanine and lactate are chemically analogous three-carbon organic acids, with alanine having an amino group where lactate has a hydroxyl group. Overall, our data confirms that that DUF4212 is required for the transporter's activity, so it is probably an additional subunit. Some members of DUF4212 are cofit with a nearby member of COG2905 (RNase T domain protein) (e.g., Sama\_1525 or Psest\_0349); we speculate that COG2905 is required for expression of the symporter.

#### **UPF0060 (PF02694, COG1742): efflux pump for thallium (I) ions**

UPF0060 is specifically important for resisting thallium (I) stress in *Cupriavidus basilensis* (RR42\_RS34240), in *Sphingomonas koreensis* (Ga0059261\_1942), and in two strains of *Pseudomonas fluorescens*. This family includes *E. coli* YnfA, which is an integral membrane protein. Nir Hus (PhD dissertation, 2005) proposed that YnfA belongs to the SMR (small multi-drug-resistance) family and SCOOP<sup>16</sup> (a family-family relationship finder) shows similarity to transporters and efflux pumps. UPF0060 is sometimes adjacent to cation efflux genes related to *czcD* or *zntA*. So, we propose that these members of UPF0060 are efflux pumps for thallium (I).

#### **UPF0126 (PF03458, COG2860): glycine transporter**

A number of representatives of UPF0126 are specifically important for utilizing glycine (PGA1\_c00920, SO1319, Sama\_2463, Psest\_1636, AO353\_13110). These genes contain a pair of UPF0126 domains and are predicted to be membrane proteins. SCOOP identified similarity between UPF0126 and TRIC channels, so we propose that this is a family of glycine transporters. This family includes *yicG* from *Escherichia coli*, which is not characterized (and which we did not identify any significant phenotypes for).



## Pathway-level Predictions:

### **DUF444 (yeaH, PF04285) and SpoVR (ycgB, PF04293): signaling with serine kinase yeaG**

As discussed in the main text, we identified strong and conserved cofitness between *yeaG*, *yeaH*, and *ycgB* in many different bacteria, including *Escherichia coli*. These genes form a single operon in most bacteria, but in *E. coli* they are broken up into two operons (*yeaGH* and *ycgB*). The biological role of YeaG is not known but it is a PrkA-like protein kinase and can phosphorylate itself or casein *in vitro*<sup>17</sup>. SpoVR is so named because a member of this family is involved in sporulation in *Bacillus subtilis*, but nothing is known about its biochemical function and the organisms that we studied do not sporulate. Since YeaG appears to be a signaling protein, we infer that YeaH and YcgB are involved in this signaling pathway as well. The mutant phenotypes of these genes are not conserved, but there is a potential commonality relating to the utilization of amino acids: these genes are detrimental for the utilization of L-arginine as the nitrogen source by *E. coli*, detrimental for the utilization of L-methionine or L-phenylalanine as the nitrogen source by *S. oneidensis* MR-1, and detrimental for the utilization of L-serine as the carbon source by *P. fluorescens* FW300-N2E3. So we suspect that the ultimate target(s) of this signaling pathway is involved in nitrogen metabolism.

### **DUF466 (PF04328, COG2879): accessory to pyruvate transport by cstA-like**

Three representatives of DUF466 were cofit with a nearby *cstA*-like protein. In *E. coli*, CstA is reported to be a peptide transporter. In the *Desulfovibrio alaskensis* G20, a *cstA*-like protein (Dde\_2007) is specifically important for fitness when pyruvate is the carbon source (data from<sup>18,19</sup>). The only strong phenotype for the DUF466 and *cstA* in our data were in *C. basilensis*, where they were specifically important for pyruvate utilization. So, we propose that DUF466 is required for pyruvate transport by a CstA-like protein.

### **DUF692 (PF05114), DUF2063 (PF09836): chlorite stress signaling proteins**

DUF692 and DUF2063 form a conserved operon that is important for chlorite resistance in *P. stutzeri* RCH2 (Psest\_0116:Psest\_0117), *Kangiella aquimarina* (B158DRAFT\_1333:B158DRAFT\_1334), and *Shewanella amazonensis* SB2B (Sama\_1305 = DUF692; DUF2063 seems to have been replaced by a hypothetical protein, Sama\_1304, which also has this phenotype). In *S. amazonensis*, DUF692 is cofit with the hypochlorite scavenging system *yedYZ* (Sama\_1893:Sama\_1892, see<sup>20,21</sup>), so DUF692 probably has another role rather than detoxifying hypochlorite directly. DUF692 is a putative xylose isomerase or epimerase-like and DUF2063 is a putative DNA binding protein. We propose that these genes are involved in sensing an aspect of chlorite stress or metal stress. Some homologs are in an operon with an extracellular type sigma factor that responds to heavy metal stress (e.g. NGO1944 from *Neisseria gonorrhoeae*, which regulates

methionine sulfoxide reductase). This suggests that DUF2063 might be an anti-anti-sigma factor rather than a DNA-binding protein. Consistent with this, in *Neisseria*, there is a putative anti-sigma factor downstream of the sigma factor (e.g., NMB2145), which is not similar to DUF2692 or DUF2063.

#### **DUF934 (PF06073): accessory protein for sulfite reduction**

This uncharacterized protein family is conserved downstream of *cysI*, the  $\beta$  subunit of sulfite reductase. In SEED, members of this family are annotated as CysX or "Oxidoreductase probably involved in sulfite reduction," but we were not able to find any published experimental evidence. They were probably annotated based on conserved gene proximity. (DUF934 does not seem to be homologous to the small ferredoxin-like CysX of *Corynebacterium glutamicum* described by Ruckert et al <sup>22</sup>; furthermore, it lacks the CxxC motifs that are conserved in CysX.) We found that DUF934 was strongly cofit with various genes in the sulfate assimilation pathway in various *Pseudomonas* (i.e., Psest\_2088) as well as in *Marinobacter adhaerens* HP15 and *C. basilensis* 4G11. In *M. adhaerens* and in *C. basilensis*, there are other sulfate assimilation genes downstream, but not in the *Pseudomonas* species, so this is not a polar effect.

#### **DUF971 (PF06155): FeS cluster maintenance protein**

In several *Shewanella* or *Pseudomonas* species, a representative of DUF971 is cofit with Mrp, BolA, and/or YggX. All of these proteins are related to FeS cluster maintenance: Mrp is an FeS loading protein; BolA is related to the Fra2 protein of *Saccharomyces cerevisiae* which is an FeS cluster protein and regulates iron levels; and YggX plays a role in oxidation resistance of FeS clusters. DUF971 is pleiotropic but a number of representatives are important for resisting cobalt (II) or paraquat; either of these might disrupt FeS clusters. Nothing is known about the biochemical function of DUF971, but in eukaryotes, DUF971 is found as a domain within gamma-butyrobetaine hydroxylase and trimethyllysine dioxygenase proteins, and it is also found within the chloroplast 4Fe4S cluster scaffold protein HCF101 (fused with an Mrp-like domain). These occurrences are consistent with DUF971 having a role in FeS cluster maintenance.

#### **DUF1302 (PF06980), DUF1329 (PF07044): export of cell wall component**

We identified cofitness for a gene cluster comprising DUF1302 (e.g., Sama\_1588 in *Shewanella amazonensis* SB2B), DUF1329 (Sama\_1589), a BNR repeat protein (COG4447; Sama\_1590), and a putative RND export protein (COG1033; Sama\_1591). The BNR repeat protein is related to a photosystem II stability/assembly factor (*ycf48* or *hcf136*) and has a beta-propeller fold (PDB 2xbg). These genes were conserved near each other and are cofit in *Kangiella aquimarina* DSM 16071 as well. In *Pseudomonas stutzeri* RCH2, the cluster is broken up into two pieces, Psest\_1122:\_1123 (DUF1302 and DUF1329) and Psest\_1923:\_1924 (BNR repeat protein and RND export protein).

Again these proteins are cofit, although in one condition (tyrosine as a carbon source), the two groups have strong and opposing phenotypes; this could imply some separation of function, or it could relate to the existence of paralogs for DUF1329 in this organism. These genes were cofit in some strains of *Pseudomonas fluorescens* as well. In both *S. amazonensis* and *P. stutzeri*, their phenotypes are pleotropic.

Several lines of evidence suggest that these proteins are in the outer membrane and affect the cell wall. DUF1329 has some similarity to the outer membrane protein LolB (e.g., the C-terminal part of Sama\_1589 is similar to CATH superfamily 2.50.20.10). In *Delftia* sp. Cs1-4, both DUF1302 and DUF1329 are reported to be components of extended outer membrane vesicles or nanopods<sup>23</sup>. BNR repeat proteins are often found in the outer membrane. And a mutant of a member of DUF1329 in *P. fluorescens* F113 is reported to have increased swimming motility<sup>24</sup>. We propose that these four proteins work together to export a component of the cell wall.

Incidentally, a number of other studies have discussed the DUF1329 family but without identifying a biochemical function or a mutant phenotype. In *Delftia* sp. Cs1-4, a DUF1329 family member (phnK) is in a genomic island for phenanthrene catabolism, where it is also near a RND efflux protein<sup>25</sup>. In *Thauera aromatica* T1, expression of a member of this family (*pipB*) is induced by p-cresol (M. Chatterjee, PhD Thesis 2012). (A DUF1302 member, *pipA*, was also induced.) In *Comamonas testosteroni*, ORF61 (DUF1329) is in a gene cluster for steroid degradation but is not required for it<sup>26</sup>. The presence of DUF1329 in these gene clusters seems to suggest a specific role in the maturation of an outer membrane transporter, or directly in the transport of aromatic compounds, but this is not what we observed in *S. amazonensis* or *P. stutzeri*.

#### **DUF1654 (PF07867): DNA repair protein**

In several strains of *Pseudomonas fluorescens*, members of DUF1654 (PF07867) are specifically important for resisting the DNA-damaging agent cisplatin (e.g., AO356\_00980). Furthermore, these genes are upstream of an endonuclease precursor, and close homologs of these proteins are predicted to be regulated by LexA, which controls the DNA damage response (e.g., PFL\_2098 from *P. fluorescens* Pf-5). These imply a role for DUF1654 in DNA repair, but we have little idea of its molecular function. Polar effects seem unlikely because DUF1654 has a strong phenotype relative to the downstream endonuclease, but we cannot rule them out entirely. DUF1654 has some similarity (SCOOP) to an anti-parallel beta barrel family (calycin\_like, PF13944).

#### **DUF2946 (PF11162): copper homeostasis protein**

In two strains of *Pseudomonas fluorescens*, DUF2946 is in an operon with and cofit with a TonB-dependent copper receptor (TIGR01778) and a Cox17-like copper chaperone (i.e., AO356\_08325 with \_08320 and \_08330; the copper chaperone is also referred to as CopZ). Some members of DUF2946 are labeled as being homologous to COG2132,

which includes *E. coli* CueO and SufI, but those are much longer proteins (around 500 amino acids, while DUF2946 is around 100 amino acids). We propose that DUF2946 is also involved in copper uptake or homeostasis.

#### **DUF3584 (PF12128): smc-like chromosome partitioning protein**

In three species of *Shewanella*, DUF3584 was specifically important for resisting cisplatin stress (i.e., Sama\_0592) and also for motility. According to HHSearch, DUF3584 is related to the N-terminal domain of Smc, or RecF or RecN, which are all involved in DNA repair or recombination. Together, this strongly suggests that DUF3584 is involved in chromosome segregation or DNA repair, but we do not have a specific proposal for its biochemical role. DUF3584 is conserved cofit with two hypothetical proteins (e.g., Sama\_0594 and Sama\_0593) that do not belong to any annotated families, and the three genes probably act together.

#### **UPF0227: hydrolase affecting the cell envelope**

UPF0227 includes the *E. coli* protein YqiA. *In vitro*, YqiA has esterase activity on palmitoyl-CoA and p-nitrophenylbutyrate<sup>27</sup>, but little is known about its biological role. Mutants in this family have diverse phenotypes, including sensitivity to bacitracin, which inhibits synthesis of a peptidoglycan precursor (*E. coli* yqiA, Sama\_3044, and Psest\_0482), and sensitivity to carbenicillin, which a  $\beta$ -lactam antibiotic and inhibits peptidoglycan synthesis (SO3900, Psest\_0482). In some *Shewanella* species (but not in *S. oneidensis* MR-1), this family is important for fitness in many defined media conditions. These phenotypes are consistent with an effect on cell wall or membrane composition. Given its esterase activity on palmitoyl-CoA, an effect on lipid composition seems most likely. A caveat is that these genes are upstream of the essential protein *parE*, and it difficult to rule out a phenotype from reducing *parE* expression (a polar effect). Although we cannot rule out this explanation, we did find that insertions within *E. coli*'s *yqiA* or within Sama\_3044 have strong phenotypes in either orientation, which makes it less likely.

## References

1. Ananthaswamy, H. N. The release of endonuclease I from *Escherichia coli* by a new cold shock procedure. *Biochem. Biophys. Res. Commun.* **76**, 289–298 (1976).
2. Nossal, N. G. & Heppel, L. A. The release of enzymes by osmotic shock from *Escherichia coli* in exponential phase. *J Biol Chem* **241**, 3055–3062 (1966).
3. Lopes, J., Gottfried, S. & Rothfield, L. Leakage of periplasmic enzymes by mutants of *Escherichia coli* and *Salmonella typhimurium*: isolation of ‘periplasmic leaky’ mutants. *J. Bacteriol.* **109**, 520–525 (1972).
4. Baba, T. *et al.* Construction of *Escherichia coli* K-12 in-frame, single-gene knockout mutants: the Keio collection. *Mol. Syst. Biol.* **2**, 2006 0008 (2006).
5. Romine, M. F., Carlson, T. S., Norbeck, A. D., McCue, L. A. & Lipton, M. S. Identification of mobile elements and pseudogenes in the *Shewanella oneidensis* MR-1 genome. *Appl Environ Microbiol* **74**, 3257–3265 (2008).
6. Deutschbauer, A. *et al.* Evidence-based annotation of gene function in *Shewanella oneidensis* MR-1 using genome-wide fitness profiling across 121 conditions. *PLoS Genet.* **7**, e1002385 (2011).
7. Rubin, B. E. *et al.* The essential gene set of a photosynthetic organism. *Proc. Natl. Acad. Sci. U.S.A.* **112**, E6634–43 (2015).
8. Kato, J.-I. & Hashimoto, M. Construction of consecutive deletions of the *Escherichia coli* chromosome. *Mol. Syst. Biol.* **3**, 132 (2007).
9. Haft, D. H. *et al.* TIGRFAMs and Genome Properties in 2013. *Nucleic Acids Res.* **41**, D387–95 (2013).
10. Romine, M. F. *et al.* Validation of *Shewanella oneidensis* MR-1 small proteins by AMT tag-based proteome analysis. *OMICS* **8**, 239–254 (2004).
11. Gimenez, R., Nuñez, M. F., Badia, J., Aguilar, J. & Baldoma, L. The gene *yjcG*, cotranscribed with the gene *acs*, encodes an acetate permease in *Escherichia coli*. *J. Bacteriol.* **185**, 6448–6455 (2003).
12. Goonesekere, N. C. W., Shipely, K. & O’Connor, K. The challenge of annotating protein sequences: The tale of eight domains of unknown function in Pfam. *Comput Biol Chem* **34**, 210–214 (2010).
13. Skerker, J. M. *et al.* Dissecting a complex chemical stress: chemogenomic profiling of plant hydrolysates. *Mol. Syst. Biol.* **9**, 674–674 (2013).
14. Krehenbrink, M., Oppermann-Sanio, F.-B. & Steinbüchel, A. Evaluation of non-cyanobacterial genome sequences for occurrence of genes encoding proteins homologous to cyanophycin synthetase and cloning of an active cyanophycin synthetase from *Acinetobacter* sp. strain DSM 587. *Arch Microbiol* **177**, 371–380 (2002).
15. Adames, K., Euting, K., Bröker, A. & Steinbüchel, A. Investigations on three genes in *Ralstonia eutropha* H16 encoding putative cyanophycin metabolizing enzymes. *Appl. Microbiol. Biotechnol.* **97**, 3579–3591 (2013).
16. Bateman, A. & Finn, R. D. SCOOP: a simple method for identification of novel protein superfamily relationships. *Bioinformatics* **23**, 809–814 (2007).
17. Tagourt, J., Landoulsi, A. & Richarme, G. Cloning, expression, purification and characterization of the stress kinase YeaG from *Escherichia coli*. *Protein Expr. Purif.* **59**, 79–85 (2008).
18. Meyer, B. *et al.* The energy-conserving electron transfer system used by *Desulfovibrio alaskensis* strain G20 during pyruvate fermentation involves reduction of endogenously formed fumarate and cytoplasmic and membrane-



- bound complexes, Hdr-Flox and Rnf. *Environ. Microbiol.* n/a–n/a (2014).  
doi:10.1111/1462-2920.12405
19. Price, M. N. *et al.* The genetic basis of energy conservation in the sulfate-reducing bacterium *Desulfovibrio alaskensis* G20. *Front Microbiol* **5**, 577 (2014).
  20. Gennaris, A. *et al.* Repairing oxidized proteins in the bacterial envelope using respiratory chain electrons. *Nature* **528**, 409–412 (2015).
  21. Melnyk, R. A. *et al.* Novel mechanism for scavenging of hypochlorite involving a periplasmic methionine-rich Peptide and methionine sulfoxide reductase. *MBio* **6**, e00233–15 (2015).
  22. Rückert, C. *et al.* Functional genomics and expression analysis of the *Corynebacterium glutamicum* fpr2-cysIXHDNYZ gene cluster involved in assimilatory sulphate reduction. *BMC Genomics* **6**, 121 (2005).
  23. Shetty, A., Chen, S., Tocheva, E. I., Jensen, G. J. & Hickey, W. J. Nanopods: a new bacterial structure and mechanism for deployment of outer membrane vesicles. *PLoS One* **6**, e20725 (2011).
  24. Navazo, A. *et al.* Three independent signalling pathways repress motility in *Pseudomonas fluorescens* F113. *Microb Biotechnol* **2**, 489–498 (2009).
  25. Hickey, W. J., Chen, S. & Zhao, J. The phn Island: A New Genomic Island Encoding Catabolism of Polynuclear Aromatic Hydrocarbons. *Front Microbiol* **3**, 125 (2012).
  26. Horinouchi, M., Kurita, T., Hayashi, T. & Kudo, T. Steroid degradation genes in *Comamonas testosteroni* TA441: Isolation of genes encoding a  $\Delta 4(5)$ -isomerase and  $3\alpha$ - and  $3\beta$ -dehydrogenases and evidence for a 100 kb steroid degradation gene hot spot. *J. Steroid Biochem. Mol. Biol.* **122**, 253–263 (2010).
  27. Kuznetsova, E. *et al.* Enzyme genomics: Application of general enzymatic screens to discover new enzymes. *FEMS Microbiol. Rev.* **29**, 263–279 (2005).