

1 **Title**

2 Defining hierarchical protein interaction networks from spectral analysis of bacterial proteomes

3

4

5

6 **Authors**

7 Mark A. Zaydman^{#,1,*}, Alexander Little⁴, Fidel Haro⁴, Valeryia Aksianiuk⁴, William J. Buchser²,
8 Aaron DiAntonio³, Jeffrey I. Gordon¹, Jeffrey Milbrandt², Arjun S. Raman^{#,4,5,*}

9

10 **Affiliations**

11 ¹Department of Pathology, Washington University School of Medicine, St. Louis, MO 63110

12 ²Department of Genetics, Washington University School of Medicine, St. Louis, MO 63110

13 ³Department of Developmental Biology, Washington University School of Medicine, St. Louis,
14 MO 63110

15 ⁴Duchossois Family Institute, University of Chicago, Chicago, IL, 60637

16 ⁵Department of Pathology, University of Chicago, Chicago, IL, 60637

17 #Contributed equally

18

19 *Correspondence to: zaydmanm@wustl.edu, araman@bsd.uchicago.edu

20

21

22 **Abstract**

23

24 Cellular phenotypes emerge from a hierarchy of molecular interactions: proteins interact to form

25 complexes, pathways, and phenotypes. We show that hierarchical networks of protein

26 interactions can be extracted from the statistical pattern of proteome variation as measured

27 across thousands of bacteria and that these hierarchies reflect the emergence of complex

28 bacterial phenotypes. We describe the mathematics underlying our statistical approach and

29 validate our results through gene-set enrichment analysis and comparison to existing

30 experimentally-derived hierarchical databases. We demonstrate the biological utility of our

31 unbiased hierarchical models by creating a model of motility in *Pseudomonas aeruginosa* and

32 using it to discover a previously unappreciated genetic effector of twitch-based motility. Overall,

33 our approach, SCALES (Spectral Correlation Analysis of Layered Evolutionary Signals),

34 predicts hierarchies of protein interaction networks describing emergent biological function using

35 only the statistical pattern of bacterial proteome variation.

36

37

38

39

40

41

42

43

44

45

46

47 **Introduction**

48
49 A fundamental problem in biology is to understand how proteins interact to create a
50 complex phenotype (Barabasi and Oltvai, 2004; Chuang *et al.*, 2010; Hartwell *et al.*, 1999;
51 Costanzo *et al.*, 2016). Biochemical and genetic studies have illustrated that complex behaviors
52 emerge from a hierarchy of protein interactions: proteins interact to form complexes, complexes
53 interact to form pathways, and pathways interact to create phenotypes (Papin *et al.*, 2004;
54 Ravasz 2009; Nurse 2008). Current strategies for identifying protein interactions span both
55 experiment and computation. Experimental methods identify protein-protein interactions (PPIs)
56 across different model systems and are continuing to evolve to be more high-throughput and
57 comprehensive (Rajagopala *et al.*, 2014; Schoenrock *et al.*, 2017; Hauser *et al.*, 2014; Koo *et*
58 *al.*, 2017; Luck *et al.*, 2020). Computational methods leverage covariation between orthologous
59 genes (orthologs) to predict PPIs (Eisen, 1998; Pellegrini *et al.*, 1999; Enrich *et al.*, 1999;
60 Valencia and Pazos, 2002). Advances in computation as well as the breadth of available data
61 have fueled continued innovation such as using statistical physical methods to infer PPIs from
62 amino-acid coevolution (Croce *et al.*, 2019; Cong *et al.*, 2019; Green *et al.*, 2021).

63 Pairwise PPIs derived from experimental and computational approaches are used to
64 infer higher-order interaction networks (Szklarczyk *et al.*, 2018). However, recent experimental
65 work has shown that protein interaction networks defined only by binary interactions are
66 incomplete, suggesting that important biological information lies in higher-order protein
67 interactions (Kuzmin *et al.*, 2018). Therefore, creating a model of genotype to phenotype
68 relationships requires the ability to identify different ‘scales’ of interactions, from pairwise to
69 higher-order, and relating these scales to describe the integration of pairwise interactions into
70 higher-order interactions. We hypothesized that (i) pairwise and higher-order information could
71 be directly extracted from the statistical pattern of covariation between orthologs and (ii) this
72 information could then be used to create a single multi-scale hierarchical model describing the
73 emergence of complex phenotypes from individual proteins.

74 We used Singular Value Decomposition (SVD) to spectrally analyze a large database of
75 non-redundant bacterial proteomes and to define a set of components of ortholog covariation
76 (an ‘SVD spectrum’). We found that covariation described by the SVD spectrum was distributed
77 according to biological scale: top components were enriched for phylogenetic information,
78 deeper components for higher-order protein interactions resembling pathways (‘indirect’
79 interactions), and deepest components for pairwise PPIs resembling physically interacting
80 protein complexes (‘direct’ interactions). Second, we introduced the concept of ‘spectral
81 correlations’, a metric representing the extent to which two bacteria or proteins share a similar
82 statistical pattern of covariation within a specific region of the SVD spectrum. We found that
83 machine-learning models trained on spectral correlation features could simultaneously predict
84 indirect and direct PPIs with higher-fidelity relative to existing computational and experimental
85 methods. Third, we introduced the concept of ‘spectral depth’—a way to relate spectral
86 correlations between different positions in the SVD spectrum. Serially thresholding spectral
87 depth defined hierarchically related protein interaction networks. We found our statistically
88 derived hierarchies reflect the emergence of complex cellular phenotypes in bacteria as
89 evidenced by gene-set enrichment analysis (GSEA) and comparison with experimentally
90 derived hierarchies in the Kyoto Encyclopedia of Genes and Genomes (KEGG). The topology of
91 these hierarchies were that bottom layers define protein interaction networks representing
92 specific functions of protein complexes, middle layers define the integration of networks within
93 bottom layers into broader functions resembling pathways, and top layers define the integration
94 of networks in middle layers into high-level functions resembling phenotypes. Finally, we
95 showed the utility of our approach by assigning global and local functions to an uncharacterized
96 protein in *Pseudomonas aeruginosa* and validating these predictions experimentally. We call
97 our approach for defining hierarchical protein interaction networks Spectral Correlation Analysis
98 of Layered Evolutionary Signals (SCALES).

99

100 **Results**

101 **Spectral decomposition of orthologous gene content among bacteria organizes**
102 **covariation from phylogenetic relationships down to pairwise PPIs**

103 Orthologs are genes in different species that originated from a common ancestral gene
104 and typically share a conserved core function. Assignment of orthologous gene groups (OGGs)
105 is a robust and computationally tractable heuristic method for inferring orthologs and has been
106 used extensively in comparative genomics (Overbeek *et al.*, 1999). To sample variation in the
107 OGG content of bacterial proteomes, we created a matrix, D^{OGG} , where each row is one of
108 7,047 UniProt bacterial reference proteomes, each column is one of 10,177 OGGs, and each
109 entry is the number of times an OGG was observed in a proteome (**Figure 1A, Table S1,**
110 **Figure 1 – figure supplement 1**, Materials and Methods) (The UniProt Consortium, 2019;
111 Huerta-Cepas *et al.*, 2017; Huerta-Cepas *et al.*, 2019). We spectrally decomposed D^{OGG} using
112 SVD (Materials and Methods) (Klema and Laub, 1980) (**Figure 1 – figure supplement 2A-C**).
113 SVD reveals patterns of correlations within the data by defining components of covariation and
114 ordering them by their ability to explain the total data variance: SVD component 1 (SVD₁)
115 explains more data-variance than SVD₂, SVD₂ more than SVD₃, etc (**Figure 1 – figure**
116 **supplement 2D**). We observed that rows of D^{OGG} corresponding to organisms sharing the
117 highest level of phylogenetic classification, i.e. phylum, clustered together when projected onto
118 SVD₁, SVD₂, SVD₃, or SVD₄ suggesting that the most dominant patterns of OGG covariation
119 arise from global phylogenetic relationships (**Figure 1 – figure supplement 3**). The vast
120 majority of the overall data variance (83%) was not explained by by SVD₁, SVD₂, SVD₃, and
121 SVD₄ taken together. We next sought to systematically interrogate what biological information, if
122 any, exists amongst deeper regions of the SVD spectrum of D^{OGG} .

123 To quantify biological information contained within different regions of the SVD
124 spectrum, we computed the mutual information (MI) shared between known biological
125 relationships spanning phylogeny to pairwise PPIs and the proximity between two proteomes or

126 proteins as defined statistically by the SVD spectrum. A high MI value indicates that the
127 statistical proximity reflects the known biological relationship, just as the clustering of proteomes
128 on SVD₁, SVD₂, SVD₃, and SVD₄ reflected phylum classification (**Figure 1 – figure supplement**
129 **3**) In the following paragraphs we detail how we defined benchmarks of known biological
130 relationships and statistical proximity between proteomes or proteins within a region of the SVD
131 spectrum.

132 Benchmarks were assembled using existing biological databases to represent a
133 hierarchy of organization from phylogenetic classification to indirect interactions in cellular
134 pathways and to direct PPIs in protein complexes (**Figure 1B**, Materials and Methods). The
135 NCBI taxonomy database was used to assemble five different phylogenetic benchmarks
136 indicating if two bacteria share the same taxonomic substrings down to the levels of phylum,
137 class, order, family or genus (**Table S2**) (NCBI Resource Coordinators, 2018). Pathway level
138 benchmarks were assembled by mining indirect PPIs found in the STRING or GO databases
139 (Szklarczyk *et al.*, 2018; The Gene Ontology Consortium, 2020). Finally, protein complex
140 benchmarks were assembled by incorporating direct PPIs identified in the Protein Databank
141 (PDB), ECOCYC database, or by analyzing amino-acid level coevolution (Coev+) (Kesler *et al.*,
142 2016; Cong *et al.*, 2019). We focused our analyses on *E. coli* K12, a well-studied model
143 organism for which a large number of high quality, experimentally supported PPIs are known
144 (**Table S3**).

145 We quantified proximity within a specific region of the SVD spectrum of D^{OGG} by
146 introducing a metric we term ‘spectral correlation’ (Materials and Methods). In detail, application
147 of SVD to D^{OGG} produces two projection matrices U^{OGG} and V^{OGG} (**Figure 1 – supplemental**
148 **figure 2A-C**). A row of U^{OGG} contains the projections of a proteome onto each SVD component:
149 column one onto SVD₁, column two onto SVD₂, and so on. Similarly, the rows of V^{OGG} contain
150 the projections of an OGG onto each SVD component. Spectral correlations are the Pearson
151 correlations between two rows of U^{OGG} or V^{OGG} . To compute spectral correlations for a specific

152 region of the SVD spectrum, Pearson correlations are computed across the columns of U^{OGG} or
153 V^{OGG} representing only the components of interest. The interpretation of a positive spectral
154 correlation is that the two proteomes or OGGs are proximal when projected onto the specified
155 set of SVD components. Because a single protein-coding gene can have multiple OGGs, we
156 approximated the projection of a protein onto each SVD component by averaging the proteins'
157 constituent OGG projections and then computed protein-protein spectral correlations as above
158 **(Figure 1 – figure supplement 4)**.

159 We computed the MI shared between each benchmark and spectral correlations within
160 all five-component windows of the SVD spectrum of D^{OGG} (Materials and Methods). To estimate
161 contributions of spurious correlations arising from finite sampling and overlap in OGG structure
162 between proteins we computed the MI for a randomized projection matrix with bootstrap support
163 **(Figure 1 – figure supplement 5)**. We observed that the MI density (bits per window) for
164 phylogenetic benchmarks declined rapidly as the spectral window was shifted deeper into the
165 SVD spectrum, quickly converging upon the MI produced by spurious correlations **(Figure 1 –**
166 **figure supplement 6A)**. In contrast, the MI density decayed more slowly with spectral position
167 for benchmarks of indirect PPIs and most slowly for benchmarks of direct PPIs. In fact, SVD₂₉₉₅₋
168 ₃₀₀₀ harbored significantly greater direct PPI MI than that produced by the model of spurious
169 correlations despite accounting for only 0.021% of data-variance ($p < 10^{-243}$ by pairwise
170 Student's T-test, **Figure 1 – figure supplement 6B)**. These results suggest that the first 3000
171 components of the SVD spectrum of D^{OGG} contain meaningful biological information.

172 We computed MI cumulative distribution functions (MI cdfs) for each benchmark across
173 the top 3000 SVD components after subtracting the contributions of spurious correlations
174 (Materials and Methods). Qualitatively, the MI cdf is the relative amount of MI gained as a
175 function of depth; reaching a value of '1' indicates that deeper regions hold no more biologically
176 meaningful information regarding a benchmark. We observed that the MI cdfs for different types
177 of benchmarks approached one in the following order: phylum, class, order, family, genus,

178 indirect PPIs, mixed indirect/direct PPIs, and direct PPIs (**Figure 1C**). Of note, the MI cdfs for
179 the three different benchmarks of direct PPIs—ECOCYC, PDB, and Coev+—were nearly
180 superimposable, suggesting that the MI estimates were robust to benchmark source. Taken
181 together, these results demonstrated that spectral decomposition of bacterial OGG variation
182 using SVD organizes covariation from phylogenetic relationships down to pairwise PPIs.

183

184 **Using spectrally resolved covariation to train random forest models to predict indirect** 185 **and direct PPIs in *E. coli* K12**

186 A well-known challenge of using covariation to infer PPIs is that protein covariation can
187 represent phylogenetic relationships, indirect interactions in pathways, direct interactions in
188 protein complexes, or noise (Schafer and Strimmer, 2005; Sul *et al.*, 2018; Nagy *et al.*, 2020).
189 The results above demonstrated that the SVD spectrum of D^{OGG} separates covariation arising
190 from these different sources. Therefore, we hypothesized that spectral correlations measured
191 across specific SVD components could be used to accurately predict PPIs and assign them to a
192 scale of organization, i.e. indirect vs direct PPI.

193 To test this hypothesis, we devised a multi-level classification task where a machine
194 learning algorithm was challenged to classify pairs of *E. coli* K12 proteins as not-interacting,
195 indirect PPI, or direct PPI (**Figure 2 – figure supplement 1**, Materials and Methods). For model
196 training and validation, a gold-standard set of well characterized *E. coli* K12 protein pairs was
197 assembled: 72,000 not-interacting pairs, 1,226 indirect PPIs, 72 direct PPIs. The indirect and
198 direct PPIs were stringently chosen based on representation in multiple databases to reduce the
199 rate of false positives in individual databases (Rajagopala *et al.*, 2014). The not-interacting pairs
200 were chosen by random selection. The relative numbers of not-interacting, indirect PPIs, and
201 direct PPIs were chosen based on prior estimates of the true proportions of these interaction
202 classes in biology (Rajagopala *et al.*, 2014; Cong *et al.*, 2019)The gold-standard dataset was
203 randomly partitioned into training (60%) and validation (40%) sets. As will be described in detail

204 below, spectral correlation features and various comparison features were extracted for the
205 protein pairs in the gold-standard dataset. For each feature, Random Forest (RF) models were
206 trained using labeled examples in the training dataset and validated using unlabeled examples
207 from the validation dataset. Additional validation tasks performed included predicting out-of-bag
208 examples in the training dataset and predicting examples in four additional comprehensive
209 benchmarks (STRING Nonbinding: $n = 14,793$ indirect PPIs, GO: $n = 79,794$ indirect or direct
210 PPIs, STRING: $n = 14,793$ indirect PPIs and $n = 5,423$ direct PPIs, PDB: $n = 809$ direct PPIs).
211 RF model performance was quantified by computing F-scores for the predictions of each
212 interaction class. F-score is the harmonic mean of precision and recall providing a holistic
213 assessment of both the accuracy and completeness of each class prediction. Finally, the entire
214 process of randomly partitioning the gold-standard dataset, training, and validation was
215 repeated 50 times to assess the reproducibility of the resultant model performance.

216 To define spectral correlation features we computed spectral correlations across three
217 sets of SVD components. The selection of these three sets was informed by the results of
218 **Figure 1C**: SVD₁₋₃₃ spanned up to the 25th percentile point of the STRING nonbinding MI cdf,
219 SVD₃₄₋₂₂₃ spanned the 25th to 75th percentile points of the STRING nonbinding MI cdf, and
220 SVD₂₂₄₋₃₀₀₀ spanned the 75th percentile point to the point at which MI estimates converged upon
221 that of the model of spurious correlations. SVD₁₋₃₃, SVD₃₄₋₂₂₃, SVD₂₂₄₋₃₀₀₀ isolated the majority of
222 the MI related to either phylogeny, indirect PPIs, or direct PPIs respectively (**Figure 2A**). For
223 each pair of proteins in the gold standard dataset, we computed spectral correlations across
224 each of these three sets of SVD components (**Figure 2 – figure supplement 2A**). We term this
225 set of three correlation values as the ‘MI windowed spectral correlations’ (MIWSCs) feature.

226 Comparison features were extracted from various existing datasets derived using
227 established methods of PPI inference. These included the experimental methods of yeast-two-
228 hybrid (‘Y2H’), affinity-purification mass spectrometry (APMS1’, ‘APMS2’) and gene epistasis
229 (‘epistasis’), as well as the computational methods of phylogeny-aware gene co-occurrence

230 ('GC'), gene neighborhood ('GN'), and gene fusion ('GF') (Szkarczyk *et al.*, 2019; Rajagopala *et*
231 *al.*, 2014; Babu *et al.*, 2014; Babu *et al.*, 2018; Hu *et al.*, 2009). As differences in the size and
232 quality of the underlying dataset can influence the fidelity of computational PPI inference, we
233 extracted additional comparison features directly from our dataset (D^{OGG}). These additional
234 features included binary MI ('b-MI'), raw covariation ('Cov'), and a Principal Components
235 Analysis (PCA) based approach considering the top k SVD components (SVD_{1-k}) (Materials and
236 Methods) (Pellegrini *et al.*, 1999; Franceschini *et al.*, 2016).

237 In our multiscale classification task, the RF models trained on MIWSCs produced
238 significantly greater F-scores across all three interaction classes compared to models trained on
239 any of the other 18 different features in all validation tasks (**Figure 2B**, **Figure 2 – figure**
240 **supplement 3**, statistical comparisons by Wilcoxon rank-sum test summarized in **Table S4**).
241 With regard to the PCA based approach, models trained on SVD₁₋₅ performed at or below the
242 median rank of all models across all three classes. F-scores for predicting indirect PPIs
243 increased as components up to SVD₁₀₀ were included (SVD₁₋₁₀₀) without improvement of the
244 prediction of direct PPIs. Including components beyond SVD₁₀₀ increased the F-scores for
245 predicting direct PPIs while decreasing F-scores for predicting indirect PPIs. Thus, the PCA
246 based approach may not be ideal because models trained on the PCA-based features poorly
247 predicted indirect PPIs, direct PPIs, or both. Taken together, these data show that the MIWSCs
248 feature is the only one of the tested features that informs high fidelity PPI predictions in *E. coli*
249 K12 across multiple scales of biological organization.

250

251 **RF models trained to predict PPIs in *E. coli* K12 using MIWSCs generalize across diverse**
252 **bacteria**

253 To test generalizability of the RF models trained on MIWSCs in *E. coli* K12 to other
254 organisms we predicted proteome-wide direct PPIs for 11 additional phylogenetically diverse
255 bacteria, including one organism (*Azotobacter vinelandii*) that was not represented in D^{OGG}

256 (Materials and Methods). For comparison, we predicted direct PPIs for the same organisms
257 using either random selection ('Random', $n = 10,000$) or by mining the high confidence
258 interactions (confidence score > 0.7) in the STRING subchannels for the methods of GC, GN, or
259 GF. To benchmark against orthogonal experimental evidence, we mined the experimental
260 subchannel of the STRING database (Materials and Methods).

261 The median precision (5th-95th percentile range in parenthesis) was significantly greater
262 for direct PPIs predicted by the RF models trained on MIWSCs: 56.6% (41.0-81.2), 0.9% (0.6-
263 5.7), 26.1% (13.3-56.9), 22.4% (18.2-34.1), or 33.6% (29.5-74.6) for MIWSC RF models,
264 random selection, GC, GN, or GF respectively (**Figure 2C**, left panel, one-sided Wilcoxon rank
265 sum test with multiple comparison found in **Table S4**). In addition, the median recall was
266 significantly greater for RF models trained on MIWSCs (**Figure 2C**, right panel, **Table S4**). The
267 recall values were low across all methods, which may reflect the previously reported high
268 number of false-positives in experimental databases (Rajagopala *et al.*, 2014; Cong *et al.*,
269 2019). Nevertheless, the MIWSC RF models predicted a median (IQR in parenthesis) of 897
270 (551 to 1609) direct PPIs per proteome.

271 In addition, we performed a head-to-head comparison for predicting direct PPIs in
272 *Mycobacterium tuberculosis* H37Rv using the MIWSC RF models versus the method of Cong
273 and colleagues that infers direct PPIs from proteome-wide amino acid coevolution (Materials
274 and Methods) (Cong *et al.*, 2019). We found that the MIWSC RF models exhibited a significantly
275 greater precision and recall when benchmarked against the STRING composite score, as done
276 previously by Cong and colleagues, or against the STRING experimental scores (**Figure 2D**,
277 Chi-squared test found in **Table S4**).

278 Taken together, these results suggest that the RF models trained to predict *E. coli* K12
279 PPIs using MIWSCs were robust and generalizable across diverse bacteria.

280

281 **Spectral decomposition of protein domain content organizes covariation according to**
282 **biological scale and informs accurate indirect and direct PPI predictions**

283 We sought to test whether the results we observed above were specific to choosing
284 OGGs as a feature of orthology. Protein domains are conserved parts of proteins that have
285 been previously used as a feature of bacterial orthology (Mistry and Finn, 2007). We defined a
286 new matrix, D^{domain} , where each row is one of the 7,047 proteomes used in the D^{OGG} matrix,
287 each column is one of 7,245 unique protein domains, and each entry is the number of times a
288 domain appears in a proteome (**Figure 3A, Table S5**).

289 The SVD spectrum derived from D^{domain} was nearly superimposable on that of D^{OGG} ,
290 suggesting that the statistical structure of covariation is similar across these different
291 orthologous features (**Figure 3B, Materials and Methods**). Similar to our analysis of D^{OGG}
292 described above, we quantified the MI shared between the various benchmarks of prior
293 biological knowledge and spectral correlations within all 5 component windows of the SVD
294 spectrum of D^{domain} . Again, we observed that the MI density (bits per window) for phylogenetic
295 benchmarks declined rapidly as the spectral window was shifted deeper into the SVD spectrum
296 (**Figure 3 – figure supplement 1A**). In contrast, the MI for indirect PPI benchmarks declined
297 more slowly and the MI for direct PPI benchmarks remained statistically significant until at least
298 SVD_{3000} (Paired Student's T-test see **Figure 3 – figure supplement 1A,B**). MI cdfs were
299 computed for each benchmark and found to mirror those derived for D^{OGG} : ordered according to
300 phylum, class, order, family, genus, indirect PPI, mixed indirect/direct PPI, direct PPI (**Figure**
301 **3C, Materials and Methods**).

302 RF models were trained to predict PPIs in *E. coli* K12 using MIWSCs computed from the
303 SVD spectrum of D^{domain} (MIWSCs_{domain}). These models were validated in the same multilevel
304 classification tasks as described above for D^{OGG} (**Figure 3 - figure supplement 2**). When
305 compared to RF models trained on features of existing computational and experimental

306 methods, the RF models trained on MIWSCs_{domain} ranking 1st, 1st, and 3rd for the classes of not-
307 interacting, indirect PPI, and direct PPI respectively (**Table S6**).

308 Taken together, these results illustrate that spectral decomposition of orthologous gene
309 content across bacterial proteomes separates covariation arising from different biological scales
310 regardless of whether orthology is defined through orthologous gene groups or protein domains.
311 As a result of this spectral separation, spectral correlations computed across sets of SVD
312 components of OGG or domain covariation can produce accurate predictions of PPIs at different
313 biological scale, i.e. indirect PPIs and direct PPIs.

314

315 **A statistically-defined hierarchy of protein interaction networks describing the emergent**
316 **phenotype of directed motility in *E. coli* K12**

317 Understanding the molecular basis of a phenotype requires (i) identifying units of
318 collective function at different biological scales and (ii) relating these scales to create a
319 hierarchical model of how a phenotype emerges from a set of proteins. A useful example is the
320 experimentally derived hierarchical model of directed motility in *E. coli* K12 (KEGG hierarchy,
321 BRITE ECO:02035). At the lowest levels in this hierarchy, physical interactions between
322 proteins create small units of collective structure and function, such as a basal body, rod, ring,
323 motor, and filament. Integration of these structures and their individual functions produces the
324 flagellum, a machine that turns to move the cell. Integration of the flagellum and the chemotaxis
325 system ultimately produces directed motility – the ability to move purposefully along a chemical
326 gradient.

327 We hypothesized that we could derive a multiscale, hierarchical model of this phenotype
328 in a purely data-driven and unbiased manner using only the SVD spectrum of D^{OGG} . To do so,
329 we first developed a model of spectral correlations between non-interacting proteins. We then
330 applied this model to identify ‘significant’ protein correlations within different regions of the SVD
331 spectrum. These significant proteins correlations represented statistically predicted protein

332 interactions. Next, we defined a metric termed 'spectral depth' as the deepest spectral position
333 to which two proteins remained significantly correlated. We posited that applying serial
334 thresholds to spectral depth would identify a tree-like hierarchy where the root of the tree is the
335 protein interaction network observed at shallower spectral depth thresholds and the branches
336 are the networks defined at deeper spectral depth thresholds. The details of creating the model
337 of non-interacting protein correlations and defining a hierarchy using spectral depth are detailed
338 below.

339 To define a model of spectral correlations between non-interacting proteins, we first
340 considered the distribution of all pairwise spectral correlations centered on SVD_{1000} for the
341 proteins encoded in the proteome of *E. coli* K12. Our rationale was that since the vast majority
342 of proteins do not interact to produce a collective function, the distribution of all-by-all spectral
343 correlations approximates that of correlations between proteins that do not functionally interact.
344 We observed that the variance of this distribution decreased rapidly as the correlation window
345 widened until about a 100 component width – motivating our choice of computing correlations
346 over sets of 100 components (**Figure 4 – figure supplement 1A,B**). We computed distributions
347 of all-by-all correlations between *E. coli* K12 proteins across windows centered on different
348 regions of the SVD spectrum and found them to be superimposable (**Figure 4 – figure**
349 **supplement 1C**). Additionally, we computed such distributions for proteins from other diverse
350 bacteria and found them to be superimposable with those derived from *E. coli* K12 (**Figure 4 –**
351 **figure supplement 1D**). These properties enabled us to define a constant threshold for
352 significant spectral correlations between two proteins across any 100 component SVD window.
353 The p-value derived from the empirical CDF of this model decreased rapidly until a threshold
354 value of 0.29 (**Figure 4 – figure supplement 1E**). Therefore, we chose the value of 0.29
355 associated with a p-value of 0.018 as the threshold of spectral correlations signifying functional
356 interactions between proteins derived from any bacterial proteome within any region of the SVD
357 spectrum.

358 To develop a hierarchical model of the directed motility phenotype in *E. coli* K12 we first
359 identified all proteins ($n = 75$) that were significantly correlated with FliC, the flagellar filament
360 protein, over the first spectral window enriched for non-phylogenetic information (SVD₃₄-SVD₁₃₄)
361 (**Figure 4 – figure supplement 2A**). For these proteins, we computed pairwise spectral
362 correlations across all 100 component windows of the SVD spectrum of D^{OGG} . We observed that
363 some pairs (e.g. MotA and MotB) remained significantly correlated as the spectral window was
364 shifted deep into the SVD spectrum while other pairs (e.g. MotA and CheR) were significantly
365 correlated only across the shallower regions of the SVD spectrum (**Figure 4 – figure**
366 **supplement 2B**). We computed the position at which the pairwise correlation first dropped
367 below the significance threshold defined by our model of correlations between non-interacting
368 proteins. We define this position as the ‘spectral depth’ of correlation. We computed the spectral
369 depth for all pairs of proteins that were significantly correlated with FliC across SVD₃₄ to SVD₁₃₄
370 (**Figure 4 – figure supplement 2C**). Apply a threshold to spectral depth generates an
371 adjacency matrix where a pixel value of ‘1’ indicates a pair of proteins that share a spectral
372 depth that is as deep or deeper than the threshold value (**Figure 4 – figure supplement 2D**).
373 This adjacency matrix can be used to construct a protein interaction network at the thresholded
374 spectral depth.

375 We constructed protein interaction networks from the adjacency matrices produced by
376 applying spectral depth thresholds of 50, 300, and 1000 (**Table S7**). At a spectral depth of 50,
377 we observed a single densely connected network devoid of obvious substructure (**Figure 4A**,
378 top panel). Gene set enrichment analysis (GSEA) indicated that this network was enriched for
379 functional terms related to ‘flagellar system’ ($p < 10^{-45}$) (Huang *et al.*, 2009; Huang *et al.*, 2009,
380 Materials and Methods). Progressing to spectral depth of 300, we observed that the network at
381 50 fractured into four discrete subnetworks (**Figure 4A**, middle panel). These subnetworks were
382 significantly enriched for terms related to ‘Chemotaxis signaling’ ($p < 10^{-15}$), ‘Flagellum’ ($p < 10^{-$
383 ⁵⁶), ‘LPS biosynthesis’ ($p < 10^{-3}$), or ‘cyclic di-GMP signaling’ ($p < 10^{-21}$). Progressing to spectral

384 depth of 1000, the subnetworks at 300 fractured further yielding 9 discrete subnetworks. Each
385 subnetwork was significantly enriched for terms related to a specific function such as ‘cyclic di-
386 GMP catabolism’ ($p < 10^{-25}$) and ‘cyclic di-GMP synthesis’ ($p < 10^{-13}$) or ‘chemotransmission’ (p
387 $< 10^{-4}$) and ‘chemoreception’ ($p < 10^{-12}$) (**Figure 4A**, bottom panel).

388 Taken together the three network diagrams derived at spectral depths 50, 300, and 1000
389 depict a hierarchy of structure and function. Subnetworks observed at deeper spectral depths
390 integrate to form the subnetworks observed at shallower spectral depths. As the subnetworks
391 coalesced, the p-value associated with GSEA remained highly significant while the ontology of
392 the significantly enriched terms changed. We interpret these observations to mean that as we
393 ascend the statistical hierarchy, molecular descriptions of new biological functions emerge from
394 the integration of functional units at lower levels.

395 We compared our hierarchical model with the model of *E. coli* K12 motility detailed within
396 the KEGG database (BR:eco02035) (Kanehisa *et al.*, 2017) (**Figure 4B**). The two models were
397 similar in several ways. First, 44 of 55 of the proteins listed in the KEGG hierarchy also
398 appeared in the statistical hierarchy. Second, 7 of the 12 categories listed in the KEGG
399 hierarchy had a one-to-one correspondence with a subnetwork of the statistical model sharing
400 an overlapping set of proteins and similar descriptive label. Finally, both hierarchies shared a
401 conserved architecture consisting of the integration of chemoreception and chemotransmission
402 into chemotaxis signaling, the integration of flagellar substructures into the flagellum, and at the
403 most global level the integration of chemotaxis and the flagellum. The most striking difference
404 was that our statistical hierarchy included subnetworks related to cyclic-di-GMP signaling and
405 LPS biosynthesis which were absent from the KEGG hierarchy. Prior experimental studies have
406 provided direct genetic evidence that these systems are involved in *E. coli* K12 motility (Paul *et*
407 *al.*, 2010, Walker *et al.*, 2004).

408 Overall, of the 75 proteins in our hierarchical model of *E. coli* K12 motility, 44 (59%) were
409 represented in the KEGG hierarchy, 28 (37%) were missing from the KEGG hierarchy but

410 supported by prior experimental evidence in the literature, and only 3 (4%) remained
411 unvalidated (CsgG, PpdD, TorS) (**Table S7**). Taken together, these results illustrate that
412 identifying the *E. coli* K12 proteins that were significantly correlated with FliC and then serially
413 thresholding their spectral depth produced a valid multiscale, hierarchical model of *E. coli* K12's
414 directed motility phenotype.

415

416 **Robustness and generalizability of defining statistical hierarchies using spectral depth**

417 We performed four additional to test the robustness and generalizability of our approach.
418 First, we characterized motility in *E. coli* K12 using MotB, the flagellar motor protein, as a query.
419 We found a similar hierarchical architecture as observed using FliC as the query with
420 chemotaxis signaling, flagellum, and cyclic-di-GMP signaling modules appearing at spectral
421 depth 300, and more fine-grained subnetworks appearing in deeper layers (**Figure 4 – figure**
422 **supplement 3, Table S8**). To test generalizability across organisms, we created a model of
423 motility in *Bacillus subtilis* 168 using its flagellar filament protein as a query (Hag). This analysis
424 again produced a hierarchical model of motility that (i) recapitulated the corresponding KEGG
425 hierarchy, (ii) identified proteins missing from the KEGG hierarchy that are known effectors of *B.*
426 *subtilis* motility, and (iii) identified a small number of putative motility effectors (**Figure 4 – figure**
427 **supplement 4, Table S9**). Next, we tested if our method could generalize to non-physically
428 coupled pathways. We produced a model of amino-acid metabolism in *E. coli* K12 using the
429 query protein HisG, an enzyme involved in Histidine biosynthesis. The resultant hierarchical
430 model identified 130 proteins that were densely connected at spectral depth 50. Progressing to
431 deeper spectral depths revealed modules corresponding to specific functions, such as amino
432 acid and nucleotide biosynthesis. At yet deeper spectral depths, modules enriched for proteins
433 involved in the synthesis of specific amino acids became evident (**Figure 4 – figure**
434 **supplement 5, Table S10**). Finally, we demonstrated that valid statistical models of *B. subtilis*
435 168 and *E. coli* K12 motility could be derived by serially thresholding spectral depth of

436 correlations within the SVD spectrum of D^{domain} (**Figure 4- figure supplement 6, Table S11,**
437 **Table S12**). Taken all together, these analyses demonstrated that serially thresholding spectral
438 depth produces a hierarchical model of biological pathways across different query proteins,
439 organisms, pathways, and orthologous features.

440

441 **Using the structure of a statistically defined hierarchy to aide in the discovery of novel**
442 **genotype-phenotype relationships**

443 The hierarchical models produced by serially thresholding spectral depth recapitulated
444 the known architecture of several well-studied biological phenotypes without incorporating any
445 prior knowledge of biological organization. This motivated the hypothesis that these models
446 could also reveal new biological organization that was not previously appreciated. We tested
447 this idea by generating a hierarchical model of motility in *Pseudomonas aeruginosa*, using it to
448 assign both a general and specific function to a previously uncharacterized protein, and
449 experimentally validating these predictions.

450 *P. aeruginosa* uses two different types of motility – propulsive motility based on a
451 flagellum and twitch motility based on a pilus (Kearns *et. al.*, 2001; Rashid and Kornberg, 2001).
452 Using PilA, a structural component of the pilus, as a query we identified a network of 141
453 proteins as significantly correlated across SVD₃₄ to SVD₁₃₄. We produced network
454 configurations for these proteins using spectral depth thresholds of 50 and 300 (**Table S13**). At
455 a spectral depth of 50 a single densely connected network was observed (**Figure 5 – figure**
456 **supplement 1A**). Significantly enriched terms for this network were ‘methyl-accepting
457 chemotaxis’ ($p < 10^{-34}$), ‘cell motility and secretion’ ($p < 10^{-33}$), ‘two-component system’ ($p < 10^{-$
458 27), ‘type IV pilus-dependent motility’ ($p < 10^{-10}$), and ‘flagellar assembly’ ($p < 10^{-6}$), suggesting
459 that these proteins are collectively involved in the global function of directed motility. At spectral
460 depth 300, the network fractured into 18 different discrete subnetworks that were enriched for
461 specific functions (**Figure 5A, Figure 5 – figure supplement 1B**). The largest subnetworks

462 were enriched for ‘methyl-accepting chemotaxis protein’ ($p < 10^{-59}$), ‘pilus motility’ ($p < 10^{-17}$), or
463 ‘bacterial flagellum’ ($p < 10^{-21}$). We noted four proteins in the pilus subnetwork (Q9I5G6,
464 Q9I5R2, Q9I0G2, Q9I0G1) that were annotated as ‘uncharacterized protein’ in UniProt (**Figure**
465 **5B**). Further review of STRING, GO, BIOCYC, and PFAM revealed only that Q9I5G6 contains a
466 ‘domain of unknown function’ (DUF4845). Based upon their membership in the pilus subnetwork
467 at spectral depth 300, we hypothesized that these proteins may contribute to the general
468 function of directed motility by affecting the specific function of pilus-based motility. Furthermore,
469 the lack of connections to the flagellum subnetwork suggested that these proteins would not
470 impact flagellar based motility.

471 To test these predictions, we screened single-gene transposon mutants of *P. aeruginosa*
472 (PAO1) for twitch-based or flagellar-based motility using established experimental assays
473 (Materials and Methods) (Kearns *et. al.*, 2001; Rashid and Kornberg, 2001, Little *et. al.*, 2018).
474 Transposon mutants of Q9I5R2, Q9I0G2, and Q9I0G1 exhibited motility that was not
475 significantly different from the parent strain in both assays (**Table S14**). In contrast, we found
476 that two different transposon mutants of Q9I5G6 exhibited significantly reduced twitch motility
477 velocity over 24, 48, and 72 hours compared to the parent strain (**Figure 5C**, $p < 10^{-4}$ by
478 Dunnett’s multiple comparisons test). This phenotype resembled that of a knockout strain of pilA
479 and was reversed upon trans-complementation. In contrast, flagellar-based motility of Q9I5G6
480 was not significantly different from that of the parent strain (**Figure 5C – inset**, $p > 0.05$). These
481 results illustrate that Q9I5G6 is a previously unappreciated effector of the global directed motility
482 phenotype in *P. aeruginosa* that specifically impacts twitch-based motility. As such, these
483 experiments provide a proof of concept of how our hierarchical models may aid in discovering
484 novel genotype-phenotype relationships.

485

486 **Discussion**

487 Connecting genotype to phenotype is a central goal in biology. Achieving this goal
488 requires understanding how the collection of proteins in a proteome interact at different scales
489 spanning protein complexes, pathways, and cellular phenotypes. Here, we have shown that a
490 hierarchy of protein interaction networks can be extracted from analyzing covariation across an
491 ensemble of bacterial proteomes. Key to this outcome were three important results. First, when
492 we spectrally decomposed proteome variation using SVD we found that biological information
493 mapped onto the SVD spectrum in a specific way: shallow components were enriched for
494 phylogenetic relationships, deeper components for functional interactions between proteins in
495 pathways, and even deeper components for physical interactions within protein complexes.
496 Second, we found that spectral correlations measured across sets of SVD components defined
497 features that informed accurate classification of protein pairs as non-interacting, indirect PPI, or
498 direct PPI. Third, we developed the concept of computing a spectral depth of correlation and
499 found that serially thresholding spectral depth produced a hierarchical model of protein
500 interaction networks. These models closely resembled the known hierarchical organization for
501 several well-studied bacterial phenotypes. Finally, we illustrated the utility of generating these
502 unbiased hierarchical models by developing a model of motility in *P. aeruginosa* and using it to
503 predict global and local functions for a previously uncharacterized protein.

504 We call our approach SCALES—Spectral Correlation Analysis of Layered Evolutionary
505 Signals. To facilitate access to SCALES and other methods described in this paper, we have
506 developed (i) a precomputed database of proteome-wide indirect (122,725,727) and direct
507 (19,546,063) PPI predictions for all 7,047 UniProt reference bacterial proteomes; (ii) a tool for
508 predicting indirect and direct PPIs for a user-input proteome; (iii) a tool for generating and
509 interrogating a hierarchical model for a query protein of interest using SCALES. All of these can
510 be found at scales.cri.uchicago.edu

511

512 *Discarding global components of covariation purifies PPI information*

513 Variation in bacterial proteomes arises from different sources of information and noise.
514 An established approach to separating information from noise is to spectrally decompose the
515 variation using SVD and then to identify which SVD components can explain more of the total
516 variance than a random process (Wigner, 1967). This leads to considering the k most global
517 SVD components (Franscescini *et. al.*, 2016). In contrast, in our study, we empirically mapped
518 the distribution of biological information across the entire SVD spectrum. Our results did not
519 match our initial expectation that the most global components would inform the highest fidelity
520 PPI predictions and minor components would solely contain noise. Instead, we found that global
521 components of covariation primarily reflected phylogeny, and PPI predictions based on these
522 components were low quality. On the other hand, we found that excluding global components
523 while retaining minor components harboring a minuscule amount of variance produced high-
524 fidelity PPI predictions. We interpret these results to mean that percent variance per component
525 does not indicate ‘importance’ of biological signal and discarding major components of
526 covariation may actually purify functionally relevant information.

527

528 *Spectral depth: a metric that extracts a hierarchy from the SVD spectrum*

529 SVD sequentially defines orthonormal vectors (components) that maximize the
530 compression of the remaining unexplained data variance. We found that there is a direct
531 mapping between the position of a component within the SVD spectrum and level in the
532 hierarchy of biological organization (**Figure 1C**). Likely this mapping reflects intrinsic differences
533 in the compressibility of biological variation arising from different hierarchical levels spanning
534 protein complexes, pathways, and phenotypes. However, SVD alone does not reveal how the
535 different levels in the hierarchy are related. Therefore, to extract a model of how protein
536 interactions are hierarchically organized to generate a phenotype, we devised the metric of
537 spectral depth, the tracking of the persistence of correlations across the SVD spectrum. We

538 found that this metric enables predicting the integration of PPIs into complex structures
539 approximating pathways and phenotypes.

540

541 Limitations

542 There are several limitations to the methods developed in this manuscript related to: 1)
543 the feature selection, 2) data requirements, and 3) lack of mechanistic insights. We discuss
544 these limitations in the following paragraphs.

545 We observed two major limitations related to the definition of an orthologous feature.
546 First, there are proteins that have no annotated conserved protein domain or OGG. For
547 example, as many as 12.3% and 16.3% of the *E. coli* K12 proteins lacked a domain or OGG
548 annotation respectively. These proteins cannot be assigned to interactions or units of function
549 by our methods. Second, there are proteins that share an infinite spectral depth of significant
550 correlation because they contain redundant annotations. These proteins may or may not
551 contribute to a shared biological function. We anticipate that continued expansion of available
552 bacterial genome sequences will improve ortholog annotation and help alleviate these
553 limitations. In addition, phylogenomics may help by providing methods that can capture
554 orthology relationships spanning both short and long timescales of evolution (Nagy *et. al.*, 2020).
555 Despite the current limitations, the sets of interactions we predicted were more accurate and
556 complete when benchmarked against existing methodologies (**Figure 2**).

557 Recently Cong and colleagues reported a method for inferring PPIs using amino-acid
558 level coevolution (Cong *et al.*, 2019). Although we did not incorporate high-resolution amino-
559 acid information into our method, we observed that our method produced more precise and
560 complete PPI predictions in *M. tuberculosis* (**Figure 2D**). At this time, it is unclear if these
561 performance differences were related to the feature of observation or methodological
562 differences. However, the results of our comparative analysis suggest that merely increasing the
563 resolution of genomic feature does not necessarily equate to more accurate PPI prediction.

564 More work will be required to determine if incorporation of amino-acid variation could help
565 resolve spectral correlations between proteins with overlapping low-resolution orthologous
566 features. Nonetheless, one important implication of using a lower resolution feature, like OGGs,
567 is our ability to compute proteome-wide PPI predictions in a matter of minutes.

568 To what degree is the ability to recover a hierarchy of biological organization dependent
569 on the ensemble of proteomes? We reason that there are two important characteristics of the
570 ensemble—the number and diversity of the proteomes. With respect to the first, we observed
571 that thousands of SVD components were required to provide a protein complex level description
572 of the 7,047 UniProt bacterial reference proteomes likely reflecting the poor compressibility of
573 this granular information. SVD can only define as many components as the smallest dimension
574 of the matrix under interrogation. In the era of genomic sequencing the number of biological
575 replicates (the rows) is typically more limiting than the number of features (the columns). As
576 such, our ability to recover protein complex level information depended on having thousands of
577 available proteomes. With respect to the second, the ensemble we used was non-redundant,
578 representing a diversity of species (**Figure 1 – figure supplement 1**). Why was this facet
579 important? Consider the extreme case where all rows represent proteome sequences from the
580 same organism. In this case, all statistical variance would be contained within a single SVD
581 component, precluding the ability to separate and relate different levels of hierarchical
582 organization. In the future, it may be useful to explore using the shape of the SVD spectrum to
583 guide sequencing efforts so that the number as well as diversity of replicates within the
584 ensemble are sufficient.

585 A notable characteristic of the methods developed here is that they are inherently
586 ‘mechanism-free’, a quality we view as both a limitation and a strength. It is a limitation in that
587 our approach can identify interactions but cannot provide insight into the function or nature of
588 those interactions. However, it is a strength because we are not limited by prior experimental

589 results and methods. As such, we believe that our approach may powerfully guide discovery of
590 novel biology.

591

592 *The potential of generalizing SCALES to other biological systems*

593 As 'big-data' in biology is becoming increasingly commonplace, to what degree are the
594 approaches developed here applicable to other biological systems? We note that the spectral
595 properties of any given dataset will be unique. As such, re-application of these methods to a
596 new dataset will require following the steps outlined in this work: creating an ensemble,
597 identifying relevant benchmarks, mapping the benchmarks onto SVD components, and
598 developing a model of random spectral correlations to define spectral depth. Beyond these
599 practical considerations, a larger question is whether the methods described here are
600 appropriate for interrogating biological systems more generally. SCALES represents a statistical
601 way to describe emergence—the integration of individual components into layers of collective
602 units of function. The property of emergence spans many biological systems, from proteins to
603 ecosystems. Thus, while it remains to be tested, it may be true that SCALES is a generally
604 useful approach to learning the hierarchical architecture of biological systems.

605

606

607 **Materials and Methods**

608

609 Generating D^{OGG}

610 All bacterial proteomes ($n = 7,047$) in the 2020_02 release of the Uniprot Reference Proteome
611 database were downloaded on 05/20/2020 (**Figure 1 – figure supplement 1**) (The UniProt
612 Consortium, 2019). OGGs were annotated using eggNOG-mapper V2 at the level of bacteria
613 ('@2') (Huerta-Cepas *et al.*, 2017; Huerta-Cepas *et al.*, 2019). An OGG count matrix was
614 assembled (D^{OGG} , **Figure 1A**) where rows were defined as proteomes, columns were defined as
615 OGGs, and the value in each cell was the number of annotations an OGG in a proteome. The
616 number of annotations was used to preserve as much information as possible versus the
617 strategy of considering binary occurrence. All OGGs present in fewer than 1% of the proteomes
618 were removed leaving 10,177 unique columns in D^{OGG} .

619

620 Singular value decomposition (SVD) of D^{OGG}

621 Singular Value Decomposition (SVD) was performed on D^{OGG} . First, the raw data matrix was
622 centered and standardized by

$$623 \quad z_{ij} = \frac{(d_{ij} - \overline{d_{[1,m],j}})}{sd_{d_{[1,m],j}}}$$

624 where z_{ij} is the ij^{th} element of the z-scored matrix Z^{OGG} , d_{ij} the ij^{th} element in the initial data
625 matrix D^{OGG} , $\overline{d_{[1,m],j}}$ and $sd_{d_{[1,m],j}}$ are the mean and standard deviation of the j^{th} column vector
626 of D^{OGG} , and m is the total number of rows in D^{OGG} ($m = 7,047$). Z^{OGG} was factorized using SVD:

$$627 \quad Z^{OGG} = U^{OGG} \Sigma^{OGG} V^{OGG T}$$

628 U^{OGG} is an $m \times K$ matrix where rows are proteomes, columns are the 'left singular vectors', and
629 each element is the projection of a proteome onto a left singular vector. Σ^{OGG} is a $m \times n$ diagonal

630 matrix where the K non-zero diagonal entries are ‘singular values’ and decrease in magnitude
631 with position along the diagonal. V^{OGG} is a $n \times K$ matrix where the rows are OGGs, columns are
632 the ‘right singular vectors’, and each element is the projection of an OGG onto a right singular
633 vector. m and n are the number of rows ($m = 7,047$) and columns ($n = 10,117$) in Z^{OGG} . K is the
634 total number of SVD components ($K = 7,047$). The fraction of data variance explained by SVD
635 component k is computed through the following equation:

636
$$var^k = \frac{\Sigma_{kk}^2}{\sum_{i=1}^K \Sigma_{ii}^2}$$

637 where var^k is the fractional variance explained by the k^{th} SVD component, Σ_{kk} and Σ_{ii} are the k^{th}
638 and i^{th} singular values respectively, and K is the total number of singular values. The plot of
639 fractional variance per component versus component number are shown for D^{OGG} in **Figure 1 –**
640 **figure supplement 2D.**

641

642 *Assembling benchmarks that collectively represent the hierarchy of biological organization*

643 The various benchmarks described within this section can be found in **Table S2, S3.**

644

645 Phylogeny benchmarks: NCBI phylogenetic strings were mapped to the NCBI taxonomic IDs
646 for each of the 7,047 bacteria represented in D^{OGG} using taxonkit 5.0
647 (<https://bioinf.shenwei.me/taxonkit/>) on 5/20/2020. Five different benchmarks were generated
648 corresponding to pairs of proteomes that share identical phylogenetic substrings down to the
649 level of phylum ($n = 5,841,696$), class ($n = 2,460,194$), order ($n = 807,338$), family ($n = 434,753$),
650 or genus ($n = 267,794$).

651

652 STRING Nonbinding benchmark: STRING database annotations were downloaded for the
653 *Escherichia coli* K12 proteome (STRING ID 511145) on 7/22/2019. A benchmark was
654 assembled to include all protein pairs (n = 14,793) with a nonzero combined STRING score that
655 did not share a 'binding' action annotation. This benchmark was expected to be enriched for
656 indirect protein-protein interactions (PPIs).

657
658 GO terms benchmark: 'Biological function' GO term annotations were mapped for the 4,391
659 proteins in the *E. coli* K12 proteome through the UniProtKB API. A benchmark was assembled
660 containing the 79,794 protein pairs that share at least 1 GO term annotation. This benchmark
661 likely contained a mixture of indirect and direct PPIs.

662
663 STRING benchmark: STRING database annotations were downloaded for the *E. coli* K12
664 proteome (STRING ID 511145) on 7/22/2019. A benchmark was assembled comprised of all (n
665 = 20,216) protein pairs with a nonzero combined STRING score. This benchmark included a
666 mixture of pairs with (n = 14,793) and without (n = 5,423) a 'binding' annotation and therefore is
667 presumed to contain a mixture of direct and indirect PPIs.

668
669 ECOCYC benchmark: A previously published benchmark included 915 pairs of *E. coli* K12
670 proteins selected from the set of complexes in the ECOCYC database after intentionally
671 excluding large complexes with greater than ten proteins to enrich for directly interacting pairs of
672 proteins (Cong *et al.*, 2019, Keseler *et al.*, 2017). This benchmark is assumed to primarily
673 represent direct PPIs.

674
675 Coev+ benchmark: A previously published set of 1,600 direct PPIs in *E. coli* K12 identified by a
676 hybrid method combining the results of amino acid coevolution (AA Coev) and prior
677 experimental data (Cong *et al.*, 2019).

678

679 PDB benchmark: A previously published set of 809 direct PPIs in *E. coli* K12 selected by the
680 criteria that they, or closely homologous proteins, have been observed to interact in a crystal
681 structure in the PDB (Cong *et al.*, 2019).

682

683 The type of biological information reflected in each benchmark is summarized in **Figure 1B**.

684

685 Computing protein-protein spectral correlations

686 A row vector in the matrix \mathbf{V}^{OGG} contains the projections of a single OGG onto each of
687 the SVD components as described by:

$$688 \quad \mathbf{v}_{i,[1,K]} = [f_i|1 \rangle \cdots f_i|k \rangle \cdots f_i|K \rangle]$$

689 where $\mathbf{v}_{i,[1,K]}$ is the i^{th} row vector of the matrix \mathbf{V}^{OGG} , f_i is the OGG represented in row i of matrix
690 \mathbf{V}^{OGG} , $f_i|k \rangle$ is the projection of f_i onto the SVD component k ($1 \leq k \leq K$), and K is the total
691 number of SVD components.

692

693 A vector representing the projections of a protein onto the SVD components was
694 generated by averaging the vectors corresponding to the projections of all OGGs annotated
695 within the protein:

$$696 \quad \boldsymbol{\Omega}_{l,[1,K]} = \frac{\sum_{f \in F} \mathbf{v}_{i_f,[1,K]}}{|F|} = [\Omega_l|1 \rangle \cdots \Omega_l|k \rangle \cdots \Omega_l|K \rangle]$$

697 where $\boldsymbol{\Omega}_{l,[1,K]}$ is the vector of averaged projections for all OGGs in protein l onto the K SVD
698 components and is used as a surrogate for the projections of protein l onto the K SVD
699 components, f is an OGG in F (the set of OGGs encoded in protein l), i_f is the index of the row in
700 the matrix \mathbf{V}^{OGG} that contains $\mathbf{v}_{i_f,[1,K]}$ (the vector of projections of OGG f onto the K SVD
701 components), $|F|$ is the number of OGGs in F , k is a single SVD component ($1 \leq k \leq K$), $\Omega_l|k \rangle$
702 is the average projection of the OGGs in protein l onto component k , and K is the total number

703 of SVD components. An example of this process is illustrated in **Figure 1 – figure supplement**
704 **4A-F**.

705 A subvector of $\Omega_{l,[a,b]}$ was extracted so as to only consider the projections of protein l
706 onto a window of SVD components as described by:

$$707 \quad \Omega_{l,[a,b]} = [\Omega_l|a \rangle \cdots \Omega_l|k \rangle \cdots \Omega_l|b \rangle]$$

708 where $\Omega_{l,[a,b]}$ is a vector of the projections of protein l onto the set of SVD components ranging
709 from component a to component b , $1 \leq a \leq K-1$, and $2 \leq b \leq K$.

710 The correlations between protein l and protein m within a spectral window was defined
711 as:

$$712 \quad \rho_{lm}^{a:b} = \text{corr}(\Omega_{l,[a,b]}, \Omega_{m,[a,b]})$$

713 where $\rho_{lm}^{a:b}$ is the correlation between proteins l and m within the spectral window ranging from
714 SVD component a to SVD component b , corr denotes the Pearson correlation, and $\Omega_{l,[a,b]}$ and
715 $\Omega_{m,[a,b]}$ are the vectors containing the projections of proteins l and m onto the SVD components
716 within the spectral window respectively. An example of this process is illustrated in **Figure 1 –**
717 **figure supplement 4G**.

718
719 A model of random spectral correlations was generated by row shuffling the matrix V^{OGG} and
720 then computing protein-protein spectral correlations as above (**Figure 1 – figure supplement**
721 **5**).

722

723 Computing Mutual Information (MI) between spectral correlations and benchmarks of biological
724 knowledge

725 For each phylogenetic benchmark, one-hundred bootstraps were generated consisting
726 of equal numbers of randomly selected pairs of proteomes that do or do not share an identical

727 phylogenetic substrating annotation in the benchmark. For each protein interaction benchmark,
728 one-thousand bootstraps were generated consisting of equal numbers of randomly selected
729 pairs of proteins that do or do not share an interaction annotation in the benchmark. For
730 bootstraps of both phylogenetic and protein interaction benchmarks, the number of pairs sharing
731 an annotation was equal to the number of pairs indicated for each respective benchmark in the
732 section 'Assembling benchmarks that collectively represent the hierarchy of biological
733 organization'.

734 Spectral correlations across all 5-component windows of the SVD spectrum between
735 component 1 and component 3000 were computed for the proteome or protein pairs in each
736 bootstrap. The uncertainty surrounding the spectral correlations within a single window was
737 calculated as the differential entropy (Cover and Thomas, 2006):

$$738 \quad H(\boldsymbol{\rho}^{a:b}) = - \sum_i \Delta p(\rho_i^{a:b}) \log_2 p(\rho_i^{a:b}) - \log_2(\Delta)$$

739 where $\boldsymbol{\rho}^{a:b}$ is the vector of spectral correlations over the window ranging from components a to
740 b for the pairs in the bootstrap, $H(\boldsymbol{\rho}^{a:b})$ is the differential entropy of $\boldsymbol{\rho}^{a:b}$, $\rho_i^{a:b}$ is a bin of
741 correlation values produced by quantizing the continuous-valued correlations in $\boldsymbol{\rho}^{a:b}$, $p(\rho_i^{a:b})$ is
742 the probability of observing a correlation value within $\rho_i^{a:b}$, and Δ is the width of the quantization
743 bins. In the present study $\Delta = 0.25$.

744 Uncertainty surrounding spectral correlations given knowledge of the phylogenetic
745 relationships or protein interactions annotated within a benchmark is described by the
746 conditional entropy:

$$747 \quad H(\boldsymbol{\rho}^{a:b} | \mathbf{c}) = p(c = 1)H(\boldsymbol{\rho}^{a:b} | c = 1) + p(c = 0)H(\boldsymbol{\rho}^{a:b} | c = 0)$$

748 where \mathbf{c} is a binary vector that assumes a value of 1 or 0 if a pair of proteomes or proteins do or
749 do not share an annotation in the corresponding benchmark respectively, $H(\boldsymbol{\rho}^{a:b} | \mathbf{c})$ is the
750 uncertainty surrounding the windowed spectral correlations given knowledge of \mathbf{c} , $p(c=1)$ and

751 $p(c=0)$ are the probability of a 1 or 0 in c respectively, and $H(\rho^{a:b}|c = 1)$ and $H(\rho^{a:b}|c = 0)$ are
752 the uncertainties surrounding the spectral correlations for subsets of pairs in the bootstrap that
753 correspond to a value of 1 or 0 in c respectively and are computed using the differential entropy
754 equation described above.

755 Finally, MI was computed as the difference between the uncertainty surrounding the
756 spectral correlations with and without knowledge of the benchmark:

$$757 \quad I(\rho^{a:b}, c) = H(\rho^{a:b}) - H(\rho^{a:b}|c)$$

758 where $I(\rho^{a:b}, c)$ is the estimate for the MI shared between the spectral correlations and the
759 benchmark. A model of random MI was generated by computing the MI shared between the
760 spectral correlations within row-shuffled versions of U^{OGG} or V^{OGG} and the benchmarks of
761 phylogeny and protein interactions respectively (**Figure 1 – figure supplement 5**). The
762 distributions of MI estimates for the different benchmarks arising from the data or random model
763 are summarized in **Figure 1 – figure supplement 6**.

764

765 Calculation of MI cumulative distribution functions (cdfs) shown in **Figure 1C**

766 Each point in the MI cdfs shown in **Figure 1C** was computed as (for the window centered on
767 component w of the SVD)

$$768 \quad cdf_w = \frac{\sum_{i=1}^w (I_i^{data} - I_i^{random})}{\sum_{i=1}^W (I_i^{data} - I_i^{random})}$$

769

770 where cdf_w is the value of the cdf at spectral position w , I_i^{data} is the MI observed in window i ,
771 I_i^{random} is the MI produced by random correlations in window i , W is the total number of windows,
772 and $1 \leq w \leq W$. Because we considered 5-component spectral windows within the first 3000
773 components, $W = 2997$.

774

775

776 *Training and validating Random Forests (RF) models for predicting PPIs in E. coli K12 using*

777 *MIWSCs*

778 *Assembling a 'gold-standard' dataset*

779 A 'gold-standard' dataset for *E. coli* K12 PPIs was assembled and consisted of 72,000 not-

780 interacting, 1,226 indirect PPIs, and 72 direct PPIs. All pairs defined as 'direct PPI' satisfied

781 three criteria: they shared amino-acid level coevolution (Coev+ benchmark), were annotated in

782 the same protein complex in the ECOCYC benchmark, and interacted in the PDB benchmark.

783 All indirect PPIs were selected based on the following criteria: they shared a 'non-binding' type

784 interaction annotation in the STRING Nonbinding benchmark, shared a 'biological function'

785 interaction in the GO benchmark, and did not share an interaction annotation in any of the

786 benchmarks of direct PPIs (Coev+, ECOCYC, or PDB). The 'not-interacting' pairs did not share

787 an interaction annotation in any of the benchmarks (GO, STRING Nonbinding, STRING, Coev+,

788 ECOCYC, or PDB). The not-interacting set was subsampled to exceed the number of physically

789 interacting pairs by 1000-fold (Rajagopala *et al.*, 2014; Cong *et al.*, 2019).

790 The gold standard pairs were randomly partitioned into training (60%) and validation

791 (40%) datasets. Fifty such random partitions were generated to assess the reproducibility of the

792 results of the machine-learning task described below.

793

794 *Training RF models*

795 RF models consisting of 100 decision trees were trained to classify pairs of proteins in *E. coli*

796 K12 as not-interacting, indirect PPIs, or direct PPIs by feeding the labeled training set examples

797 to the TreeBagger algorithm (Matlab, v2020a). This process was repeated for each random

798 partition of the gold-standard dataset yielding an ensemble of 50 RF models per feature.

799 *Validating RF models using the validation dataset*

800 Each trained RF model was subjected to three validation tasks of classifying interaction types
801 for unlabeled pairs of *E. coli* K12 proteins in the validation portion of the gold-standard (40%)
802 **(Figure 2 – figure supplement 1, Figure 2 – figure supplement 2A)**. The model performance
803 was evaluated by computing an F-score for each interaction type (not-interacting, indirect PPIs,
804 direct PPIs), where F-score is the harmonic mean of precision and recall, precision is the ratio of
805 the number of correctly predicted interactions within a class to the total number of predicted
806 interactions in a class, and recall is the ratio of the number of correctly predicted interactions
807 within a class to the total number of interactions of that class.

808

809 *Training and validating RF models on quantitative features of existing methods*

810 For each feature extracted from existing methods described below, RF models were trained and
811 validated using the identical protocol as for MIWSCs (described in the section *Training and*
812 *validating Random Forests (RF) models for predicting PPIs in E. coli K12 using MIWSCs*).

813

814 *Existing experimental features*

815 Previously published datasets derived from large scale experimental PPI screens in *E.*
816 *coli* K12 were used to generate a set of four different experimental features including: gene
817 interaction scores from a gene epistasis screen (Epistasis, n = 41,820), sum log-likelihood
818 scores from an affinity purification mass spectrometry screen (APMS1, n = 12,801), protein
819 interaction scores from an affinity purification mass spectrometry screen (APMS2, n = 291), and
820 binary pairs from a yeast-two hybrid screen (Y2H, n=1,766) (Rajagopala *et al.*, 2014; Babu *et*
821 *al.*, 2014; Babu *et al.*, 2018; Hu *et al.*, 2009).

822

823 *Existing computational features*

824 Gene co-occurrence, gene fusion, and gene neighborhood subscores for *E. coli* K12
825 (STRING ID 511145) were extracted from the STRING database on 7/22/2019 (Szkłarczyk *et*
826 *al.*, 2019; Rajagopala *et al.*, 2014; Babu *et al.*, 2014; Babu *et al.*, 2018; Hu *et al.*, 2009). Any
827 pairs without an interaction annotation in the STRING database were assigned a subscore of
828 zero.

829
830 *Binary MI (b-MI) feature*

831 The b-MI feature was modeled after the popular phylogenetic profiling method of
832 Pelligrini and colleagues (Pellegrini *et al.*, 1999). First, a binary OGG content matrix was defined
833 as follows:

$$834 \quad \mathbf{B}^{OGG} = \begin{cases} 1, & \mathbf{D}^{OGG} > 0. \\ 0, & \text{otherwise.} \end{cases}$$

835 Where \mathbf{B}^{OGG} is the binary OGG content matrix and has the same dimensions as \mathbf{D}^{OGG} .

836
837 The phylogenetic profile of an OGG was defined as:

$$838 \quad \mathbf{pp}_j = \mathbf{B}^{OGG}_{[1,m],j}$$

839 where \mathbf{pp}_j is the phylogenetic profile of OGG j , and $\mathbf{B}^{OGG}_{[1,m],j}$ and m are the j^{th} column vector
840 and number of rows in the \mathbf{B}^{OGG} respectively. For each pair of proteins in the proteome of *E. coli*
841 K12, a phylogenetic profiling feature of protein coevolution was defined as the average of the MI
842 shared between the profiles of all pairs of OGGs encoded in the two proteins:

$$843 \quad b - MI_{lp} = \frac{\sum_{j \in J} \sum_{k \in K} I(\mathbf{pp}_j, \mathbf{pp}_k)}{|J| * |K|}$$

844 Where $b - MI_{lp}$ is the MI shared between the phylogenetic profiles of protein l and protein p , J
845 and K are the sets of OGGs encoded in proteins l and p respectively, j and k are elements of J
846 and K respectively, \mathbf{pp}_j and \mathbf{pp}_k are the phylogenetic profiles of j and k respectively,
847 $I(\mathbf{pp}_j, \mathbf{pp}_k)$ is the MI shared between \mathbf{pp}_j and \mathbf{pp}_k computed using Shannon's classic

848 formulation for the MI between two discrete random variables (Shannon, 1970; Cover and
849 Thomas, 2006), and $|J|$ and $|K|$ are the number of elements in J and K respectively.

850

851 *Covariation feature*

852 The covariation between a pair of OGGs was described by:

853
$$Cov_{jk} = \frac{1}{m} \sum_{i=1}^m (d_{ij} - \overline{d_{i,[1,n]}})(d_{ik} - \overline{d_{i,[1,n]}})$$

854 where Cov_{jk} is the covariation between OGGs j and k , m and n are the number of proteomes
855 (rows) and OGGs (columns) in D^{OGG} respectively, d_{ij} and d_{ik} are the number of annotations of
856 OGGs j and k in proteome i respectively, and $\overline{d_{i,[1,n]}}$ is the average number of OGG annotations
857 in proteome i obtained by averaging the corresponding row vector in D^{OGG} . For each pair of
858 proteins in the proteome of *E. coli* K12, protein covariation was defined as the average of the
859 covariation shared between all pairs of OGGs encoded in the two proteins:

860
$$Cov_{lp}^{protein} = \frac{\sum_{j \in J} \sum_{k \in K} Cov_{jk}}{|J| * |K|}$$

861 where $Cov_{lp}^{protein}$ is the covariation feature of interaction between protein l and protein p , J and
862 K are the sets OGGs encoded in proteins l and p respectively, j and k are elements of J and K
863 respectively, Cov_{jk} is the covariation between OGGs j and k , and $|J|$ and $|K|$ are the number of
864 elements in J and K respectively.

865

866 *PCA-based spectral correlations features*

867 These features were inspired by the approach of Franceschini and colleagues and the
868 typical use of SVD to produce a low rank approximation of the initial data matrix in an effort to
869 'denoise' the data (Franceschini *et al.*, 2016). For each pair of proteins in the *E. coli* K12
870 proteome spectral correlations were computed as described in the section 'Computing protein-

871 *protein spectral correlations*' over windows ranging from component 1 to component k , where k
872 = 5, 10, 20, 50, 100, 500, 1000, 5000, or 7047.

873

874

875 *Validating RF models in two additional validation tasks*

876 *Training dataset task*

877 Each decision tree within an RF model was tasked with predicting interaction classes for
878 the out-of-bag examples from the training datasets. F-scores were computed for the consensus
879 predictions of each model.

880

881 *Comprehensive benchmark task*

882 Biological interactions are typically sparse: the number of not-interacting pairs of proteins
883 vastly outnumber the number of interacting pairs. As such, we desired to challenge each of the
884 RF models in a validation task reflective of this asymmetry. To do so, each RF model was
885 tasked with predicting classes for all pairs of proteins in the *E. coli* K12 proteome after exclusion
886 of pairs used in the gold-standard dataset. These predictions were validated against four
887 different comprehensive benchmarks: the indirect PPIs in the STRING Nonbinding benchmark
888 ($n = 5,423$ indirect PPIs, 9,637,213 not-interacting), the mixed indirect/direct PPIs in the GO ($n =$
889 79,794 indirect or direct PPIs, 9,562,842 not-interacting) and STRING benchmarks ($n = 20,216$
890 indirect or direct PPIs, 9,622,420 not-interacting), and the direct PPIs in the entire PDB
891 benchmark ($n = 809$ direct PPIs, 9,614,827 not-interacting).

892

893 *Predicting proteome-wide direct PPIs for 11 phylogenetically unrelated bacteria*

894 *Proteomes represented in D^{OGG}*

895 Each of the fifty RF models trained to classify interactions in *E. coli* K12 using MIWSCs
896 were used to predict proteome-wide indirect and direct PPIs in the following bacteria (Uniprot
897 Proteome ID, NCBI taxonomy ID in parentheses): *Aliivibrio fischeri* ES114 (UP000000537,
898 312309), *Azotobacter vinelandii* DJ (UP000002424, 322710), *Bacillus subtilis* 168
899 (UP000001570, 224308), *Caulobacter vibrioides* (UP000053705, 155892), *Helicobacter pylori*
900 26695 (UP000000429, 85962), *Mycobacterium tuberculosis* H37Rv (UP000001584, 83332),
901 *Mycoplasma genitalium* G37 (UP000000807, 243273), *Pseudomonas fluorescens* F113
902 (UP000005437, 1114970), *Staphylococcus aureus* NCTC 8325 (UP000008816, 93061),
903 *Streptomyces coelicolor* A3(2) (UP000001973, 100226), *Synechocystis sp.* PCC 6803
904 (UP000001425, 1111708). For each proteome, a set of consensus PPIs was defined as those
905 for which a majority of the models (> 25) produced the same classification of 'indirect PPI' or
906 'direct PPI'.

907

908 *Proteomes not represented in D^{OGG}*

909 To predict interactions for a proteome that was not represented in D^{OGG} (ex. *Azotobacter*
910 *vinelandii* DJ, UP000002424, 322710), OGGs were mapped using EggNOG mapper V2 and
911 MIWSCs were extracted using the OGG projections in V^{OGG} (Huerta-Cepas, 2017; Huerta-
912 Cepas, 2019). These features were used to predict proteome-wide indirect and direct PPIs as
913 described for the Uniprot Reference Proteomes above.

914

915 *Validating direct PPI predictions against experimental evidence in the STRING database*

916 The predicted direct PPIs were benchmarked against the sets of interactions in the
917 STRING database with a non-zero experimental subchannel score for *E. coli* K12 and the
918 eleven additional organisms described above.

919

920 *A head-to-head comparison with the approach of Cong and colleagues*

921 Cong and colleagues have provided proteome-wide PPI predictions for *E. coli* K12 and
922 *Mycobacterium tuberculosis* H37Rv (Cong *et al.*, 2019). Their predictions of *E. coli* PPIs were
923 based on amino acid coevolution supplemented with existing knowledge ('Coev+'). In contrast,
924 their predictions of PPIs in *M. tuberculosis* were based on amino acid coevolution alone
925 ('Coev'). Therefore, for a head-to-head comparison, we compared the predictions produced by
926 our RF models trained on MIWSCs with their PPI predictions in *M. tuberculosis*. We
927 benchmarked these interactions using two strategies. The first strategy mirrored that used by
928 Cong and colleagues, computing the fraction of interactions assigned a STRING combined
929 score of 0, 0-0.4, or > 0.4. The second strategy used orthogonal experimental evidence by
930 computing the fraction of interactions assigned a STRING experimental subchannel score of 0
931 and > 0.

932
933 Generating D^{domain}

934 PFAM domain annotations were downloaded from the UniProt database on 05/12/2020
935 (The UniProt Consortium, 2019). A domain count matrix was assembled (D^{domain} , **Figure 3A**)
936 where rows were defined as proteomes ($n = 7,047$, the same proteomes as described in **Figure**
937 **1 – figure supplement 1** and used to create D^{OGG}), columns were defined as domains, and the
938 value in each cell was the number of annotations of a domain in a proteome. All domains
939 present in fewer than 1% of the proteomes were removed leaving 7,245 unique columns in
940 D^{domain} .

941
942 Performing SVD on D^{domain}

943 SVD was performed on D^{domain} following the same protocol described in the section
944 'Singular value decomposition (SVD) of D^{OGG} '. The SVD spectrum of D^{domain} is displayed in
945 **Figure 3B**, overlaid on the SVD spectrum derived from D^{OGG} .

946

947 Relating the structure of protein covariation defined by domain-based features with the
948 hierarchy of biological organization as shown in **Figure 3C**

949

950 *Computing domain-based proteome-proteome and protein-protein spectral correlations*

951 Domain-based proteome-proteome and protein-protein spectral correlations were
952 computed in an identical fashion as described for OGG-based spectral correlations as described
953 in the sections ‘*Computing proteome-proteome spectral correlations*’ and ‘*Computing protein-*
954 *protein spectral correlations*’ respectively.

955

956 *Computing MI shared between domain-based spectral correlations and benchmarks of*
957 *biological organization*

958 The MI shared between spectral correlations within all 5-component windows of the SVD
959 spectrum of D^{domain} , ranging from component 1 to component 3000, and benchmarks of
960 biological organization was computed in the identical way and for the identical benchmarks as
961 described in the section *Computing Mutual Information (MI) between spectral correlations and*
962 *benchmarks of biological knowledge* and is shown and compared to a random model in **Figure**
963 **3 – figure supplement 1**.

964

965 *Calculation of domain-based MI cdfs shown in **Figure 3C***

966 MI cdfs shown in **Figure 3C** were computed in an identical fashion as described in the
967 section ‘*Calculation of MI cumulative distribution functions (cdfs) shown in **Figure 1C***’.

968

969 *Training and validating MIWSC_{domain}-based RF models to infer PPIs*

970 RF models were trained using MIWSC_{domains}, validated, and compared to existing methods
971 according to the same protocols described in the sections *Training and validating Random*
972 *Forests (RF) models for predicting PPIs in E. coli K12 using MIWSCs, Training and validating*

973 *RF models on quantitative features of existing methods, and Validating RF models in two*
974 *additional validation tasks.*

975

976 *Gene-set enrichment analysis performed on statistical model of E. coli K12 motility*

977 Gene-set enrichment analysis (GSEA) was performed on the sets of proteins defined by
978 the statistical modules using DAVID analysis (v6.8). The ontological term with the lowest p-
979 value is indicated for each statistical module shown in **Figure 4**. A full list of significant
980 ontological terms and their associated p-values for each statistical module is listed in **Table S4**.

981

982 *Evaluating the robustness and generalizability of predicting the hierarchical organization of*
983 *biological pathways using spectral correlations*

984 For the additional analyses characterizing motility in *E. coli* K12 using MotB, characterizing
985 motility in *B. subtilis* using Hag, characterizing amino acid metabolism in *E. coli* K12 using HisG,
986 and characterizing motility in *E. coli* K12 and *B. subtilis* using FliC and Hag respectively with
987 domain-based spectral correlations, the network graphs shown in **Figure 4 – figure**
988 **supplement 3-6** were generated following the identical approach described in the section
989 '*Predicting the hierarchical organization of the motility pathway in E. coli K12*'. Note that for the
990 domain-based characterization, a separate null model for spectral correlations in D^{domain} was
991 developed and applied as described in the section *Developing a null model of random protein-*
992 *protein spectral correlations within the SVD spectrum of D^{OGG}* . Gene-set enrichment analysis
993 was performed for each example of a statistically derived pathway in an identical fashion as
994 described in the section *Gene-set enrichment analysis performed on statistical model of E. coli*
995 *K12 motility*.

996

997 Assaying strains of *P. aeruginosa* for pilus and flagellar motility

998 All *P. aeruginosa* strains used in this study were ordered from the Manoil Lab. All strains
999 were grown at 37°C on LB supplemented with 25µg/ml irgasan and gentamicin (75 µg/ml) as
1000 necessary. *E. coli* XL1-Blue was maintained on LB agar plates with gentamicin (15 µg/ml) as
1001 necessary.

1002 *P. aeruginosa* transposon mutants of interest were ordered from the Manoil Lab. *P.*
1003 *aeruginosa* growth was at 37°C on LB supplemented with 25µg/ml irgasan and gentamicin (75
1004 µg/ml) as necessary. Strains were assayed for subsurface twitching motility as previously
1005 described (Alm and Mattick, 1995; Little *et al*, 2018). Strains were grown overnight and stab
1006 inoculated in the interstitial space between the basal surface of 1.0% LB agar and a plastic petri
1007 dish. Plates were incubated for 48 hours at 37°C. Agar was removed and cells attached to the
1008 plate were stained with 0.5% crystal violet; twitch zone diameter was measured and plates were
1009 imaged.

1010 Surface twitching motility assays were performed as previously described (Little *et al.*,
1011 2018; Kearns *et al.*, 2001). *P. aeruginosa* strains of interest were grown overnight and
1012 concentrated in morpholinepropanesulfonic acid (MOPS) buffer (10mM MOPS, 8mM MgSO₄,
1013 pH 7.6). A 2.5µl volume of the MOPS buffered bacterial suspension was spotted onto buffered
1014 twitching motility plates (10mM Tris, 8mM MgSO₄, 1mM NaPO₄, 0.5% glucose, 1.5% agar, pH
1015 7.6) and was incubated for 24 hours at 37°C. The twitching zone was measured and imaged.

1016 Swimming motility was performed as previously described (Rashid and Kornberg, 2000).
1017 Overnight cultures were stab inoculated into the surface of LB-0.3% bacto agar and were
1018 incubated for 24 hours at 37°C. The resulting swimming zone was measured.

1019 For complementation of genes of interest into *P. aeruginosa* strains, the
1020 complementation vector pBBR1-MCS5-PA0769 was created using Gibson assembly. The

1021 vector was transferred to *P. aeruginosa* by electroporation using 2.2kV in a 2mm gap cuvette,
1022 and subsequent selection using gentamicin.

1023

1024

1025

1026

1027

1028 **Tables with titles and Legends**

1029

1030 **Table S1: D^{OGG} matrix, related to Figure 1A.**

1031

1032 **Table S2: NCBI taxonomic strings for each organism used to generate phylogenetic**
1033 **benchmarks, related to Figure 1B.**

1034

1035 **Table S3: Benchmarks of PPIs in *E. coli* K12, related to Figure 1B.**

1036

1037 **Table S4: Data and statistical support for RF model validation studies, related to Figure 2,**
1038 **Figure 2—Figure Supplement 3.**

1039

1040 **Table S5: D^{domain} matrix, related to Figure 3.**

1041

1042 **Table S6: Data and statistical support for domain-based RF model validation studies,**
1043 **related to Figure 3—Figure Supplement 2.**

1044

1045 **Table S7: Data pertaining to Figure 4.**

1046

1047 **Table S8: Data related to Figure 4—Figure Supplement 3.**

1048

1049 **Table S9: Data related to Figure 4—Figure Supplement 4.**

1050

1051 **Table S10: Data related to Figure 4—Figure Supplement 5.**

1052

1053 **Table S11: Data related to Figure 4—Figure Supplement 6.**

1054

1055 **Table S12: Data related to Figure 4—Figure Supplement 6.**

1056

1057 **Table S13: Data related to Figure 5A,B and Figure 5—Figure Supplement 1.**

1058

1059 **Table S14: Data related to Figure 5C.**

1060

1061

1062

1063 **All Tables can be downloaded at github.com/arjunsraman/Zaydman_et_al as .zip files**
1064 **(Tables_S1_to_S4.zip, Tables_S5_to_S14.zip).**

1065

1066

1067 **Acknowledgments**

1068 We thank Robert Y. Chen, Adam Bailey, Nima Mosammaparast, and Jacqueline Payton for
1069 substantial discussion regarding this manuscript. We thank Rama Ranganathan for a critical
1070 reading of the manuscript as well as in-depth discussion. We thank Dinanath Sulakhe (Center
1071 for Research Informatics (CRI), University of Chicago) for assisting in producing the web
1072 application tools described in this manuscript. We thank Sam Light, Sampriti Mukherjee, and
1073 Eric Pamer for helpful discussions regarding experiments performed.

1074

1075 **Author Contributions**

1076 M.A.Z. and A.S.R. conceived the project, developed the mathematical approaches described,
1077 wrote the code for conducting analyses, performed data analysis, and assembled the
1078 manuscript. A.L., performed the experiments regarding *P. aeruginosa*. W.J.B. provided technical
1079 expertise and wrote code for annotation and parsing and contributed to the writing of the
1080 manuscript. A.D., J.I.G., and J.M. contributed to this manuscript by providing critical feedback,
1081 in-depth discussion, and contributing to the writing of this manuscript.

1082

1083 **Competing Interests**

1084 None relevant to this manuscript.

1085

1086 **References**

1087

1088 Barabasi AL, Zoltan O, (2004) **Network Biology: Understanding the cell's functional**
1089 **organization** *Nat. Rev. Gen.* **5**:101-113. <https://doi.org/10.1038/nrg1272>

1090

1091 Chuang HY, Hofree M, Ideker TA, (2010) **A decade of systems biology.** *Annu. Rev. Cell*
1092 *Dev. Biol.* **26**:721-744. <https://doi.org/10.1146/annurev-cellbio-100109-104122>

1093

1094 Hartwell LH, Hopfield JJ, Leibler S, Murray AW, (1999) **From molecular to modular cell**
1095 **biology.** *Nature* **402**:C47-C52. <https://doi.org/10.1038/3501154>

1096

1097 Costanzo M. *et al*, (2016) **A global genetic interaction network maps a wiring diagram**
1098 **of cellular function.** *Science* **353**:aaf1420. <https://doi.org/10.1126/science.aaf1420>

1099

1100 Papin J, Reed J, Paulsson BO, (2004) **Hierarchical thinking in network biology: the**
1101 **unbiased modularization of biochemical networks.** *Trends in Biochemical Sciences*
1102 **29**:641-647. <https://doi.org/10.1016/j.tibs.2004.10.001>

1103

1104 Ravasz E, (2009) **Detecting Hierarchical Modularity in Biological Networks.** *Methods*
1105 *Mol Biol.* **54**:145-60. https://doi.org/10.1007/978-1-59745-243-4_7

1106

1107 Nurse P, (2008) **Life, logic and information.** *Nature* **454**:424-426.
1108 <https://doi.org/10.1038/454424a>

1109

1110 Rajagopala S *et al.*, (2014) **The binary protein-protein interaction landscape of**
1111 ***Escherichia coli*.** *Nat. Biotechnol.* **32**:285-290. <https://doi.org/10.1038/nbt.2831>

1112

1113 Scheonrock A *et al.*, (2017) **Evolution of protein-protein interaction networks in yeast.**
1114 *PLoS One* **12**:e0171920. <https://doi.org/10.1371/journal.pone.0171920>

1115

1116 Hauser R *et al.*, (2014) **A second-generation protein-protein interaction network of**
1117 ***Helicobacter pylori*.** *Mol. Cell. Proteomics* **13**:1318-1329.
1118 <https://doi.org/10.1074/mcp.O113.033571>

1119

1120 Koo B *et al.*, (2017) **Construction and analysis of two genome-scale deletion libraries**
1121 **for *Bacillus subtilis*.** *Cell Syst.* **4**:291-305.e.7. <https://doi.org/10.1016/j.cels.2016.12.013>

1122

1123 Luck K *et al.*, (2020) **A reference map of the human binary protein interactome.** *Nature*
1124 **580**:402-408. <https://doi.org/10.1038/s41586-020-2188-x>

1125

1126 Eisen J, (1998) **Phylgenomics: Improving functional predictions for uncharacterized**
1127 **genes by evolutionary analysis.** *Genome Research* **8**:163-167.
1128 <https://doi.org/10.1101/gr.8.3.163>

1129

1130 Pellegrini M *et al.*, (1999) **Assigning protein functions by comparative genome**
1131 **analysis: Protein phylogenetic profiles.** *Proc. Natl. Acad. Sci.* **96**:4285-
1132 4288. <https://doi.org/10.1073/pnas.96.8.4285>

1133

1134 Enright A, Iliopoulos I, Kyrpides N, Ouzounis C, (1999) **Protein interaction maps for**
1135 **complete genomes based on gene fusion events.** *Nature* **402**:86-90.

1136 <https://doi.org/10.1038/47056>

- 1137
1138
1139 Valencia A, Pazos F, (2002) **Computational methods for the prediction of protein**
1140 **interactions.** *Curr. Opinion in Struct. Biol.* **12**: 368-373. [https://doi.org/10.1016/s0959-](https://doi.org/10.1016/s0959-440x(02)00333-0)
1141 [440x\(02\)00333-0](https://doi.org/10.1016/s0959-440x(02)00333-0)
- 1142 Croce G, Gueudre T, Cuevas M, Keidel V, Figliuzzi M, Szurmant H, Weigt M, (2019) **A multi-**
1143 **scale coevolutionary approach to predict interactions between protein domains.** *PLoS*
1144 *Comput. Biol.* **15**:e1006891. <https://doi.org/10.1371/journal.pcbi.1006891>
- 1145 Cong Q, Anishchenko I, Ovchinnikov S, Baker D, (2019) **Protein interaction networks**
1146 **revealed by proteome coevolution.** *Science* **365**:185-189.
1147 <https://doi.org/10.1126/science.aaw6718>
1148
- 1149 Green A, Elhabashy H, Brock K, Maddamsetti R, Kohlbacher O, Marks D, (2021) **Large-**
1150 **scale discovery of protein interactions at residue resolution using co-evolution**
1151 **calculated from genomic sequences.** *Nat. Commun.* **12**:1396.
1152 <https://doi.org/10.1093/bioinformatics/bty862>
1153
- 1154 Szklarczyk D *et al.*, (2018) **STRING v11: protein-protein association networks with**
1155 **increased coverage, supporting functional discovery in genome-wide experimental**
1156 **datasets.** *Nucleic Acids Res.* **47**:D607-D613. <https://doi.org/10.1093/nar/gky1131>
1157
- 1158 Kuzmin E *et al.*, (2018) **Systematic analysis of complex genetic interactions.** *Science*
1159 **360**:eaao1729. <https://doi.org/10.1126/science.aao1729>
1160
- 1161 Kanehisa M, Goto S, (2000) **KEGG: Kyoto Encyclopedia of Genes and Genomes.**
1162 *Nucleic Acids Res.* **28**:27-30. <https://doi.org/10.1093/nar/28.1.27>
1163
- 1164 Kanehisa M, (2019) **Toward understanding the origin and evolution of cellular**
1165 **organisms.** *Protein Sci.* **28**:1947-1951. <https://doi.org/10.1002/pro.3715>
1166
- 1167 Kanehisa M, Furumichi M, Sato Y, Ishiguro-Watanabe M, Tanabe M, (2021) **KEGG:**
1168 **integrating viruses and cellular organisms.** *Nucleic Acids Res.* **49**:D545-D551.
1169 <https://doi.org/10.1093/nar/gkaa970>.
1170
- 1171 Overbeek R, Fonstein M, D'Souza M, Pusch G, Maltsev N, (1999) **The use of gene**
1172 **clusters to infer functional coupling.** *Proc. Nat'l. Acad. Sci.* **96**:2896-2901.
1173 <https://doi.org/10.1073/pnas.96.6.2896>
1174
- 1175 The UniProt Consortium, (2019) **UniProt: a worldwide hub of protein knowledge.** *Nucleic*
1176 *Acids Res.* **47**:D506-515. <https://doi.org/10.1093/nar/gky1049>
1177
- 1178 Letunic I, Bork P, (2019) **Interactive Tree of Life (iTOL) v4: recent updates and new**
1179 **developments.** *Nucleic Acids Res.* **47**:W256-W259. <https://doi.org/10.1093/nar/gkz239>
1180
- 1181 Huerta-Cepas J, Forslund K, Coelho LP, Szklarczyk D, Jensen L, Mering C, Bork P, (2017)
1182 **Fast Genome-Wide Functional Annotation through Orthology Assignment by**
1183 **eggNOG-Mapper.** *Mol. Biol. Evol.* **34**:2115-2122. <https://doi.org/10.1093/molbev/msx148>
1184
1185

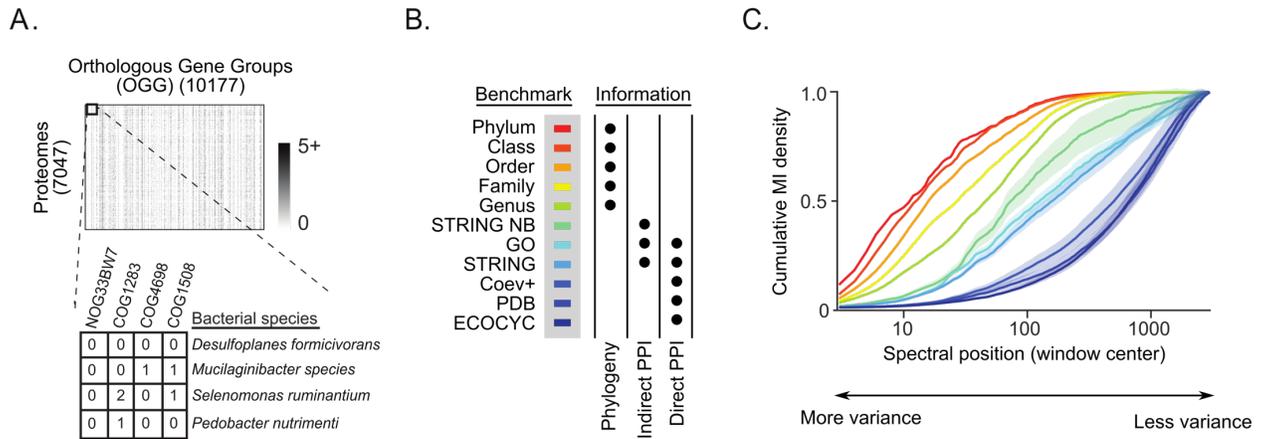
- 1186 Huerta-Cepas J, *et al.* (2019) **eggNOG 5.0: a hierarchical, functionally and**
1187 **phylogenetically annotated orthology resource based on 5090 organisms and 2502**
1188 **viruses.** *Nucleic Acids Res.* **47**:D309-D314. <https://doi.org/10.1093/nar/gky1085>
1189
- 1190 Klema V, Laub A, (1980) **The singular value decomposition: Its computation and some**
1191 **applications.** *IEEE Transactions on Automatic Control* **25**:164-176.
1192 <https://doi.org/10.1109/TAC.1980.1102314>
1193
- 1194 NCBI Resource Coordinators (2018) **Database resources of the National Center for**
1195 **Biotechnology Information.** *Nucleic Acids Res.* **46**:D8-D13.
1196 <https://doi.org/10.1093/nar/gkv1290>
1197
- 1198 The Gene Ontology Consortium (2020) **The Gene Ontology resource: enriching a Gold**
1199 **mine.** *Nucleic Acids Res.* **49**:D325-334. <https://doi.org/10.1093/nar/gkaa1113>
1200
- 1201 Kesler I, *et al.* (2016) **The EcoCyc database: reflecting new knowledge about**
1202 ***Escherichia coli* K-12.** *Nucleic Acids Res.* **45**:D543-550.
1203 <https://doi.org/10.1093/nar/gkw1003>
1204
- 1205 Cong Q, Anishchenko I, Ovchinnikov S, Baker D, (2019) **Protein interaction networks**
1206 **revealed by proteome coevolution.** *Science* **365**:185-189.
1207 <https://doi.org/10.1126/science.aaw6718>
1208
- 1209 Schafer, J., and Strimmer, K. (2005). **An empirical Bayes approach to inferring large-**
1210 **scale gene association networks.** *Bioinformatics* **6**, 185-189.
1211 <https://doi.org/10.1093/bioinformatics/bti062>
1212
- 1213 Sul, J.H., Martin, L.S., and Eskin, E. (2018). **Population structure in genetic studies:**
1214 **confounding factors and mixed models.** *PLoS. Genetics*,
1215 <http://doi.org/10.1371/journal.pgen.1007309>
1216
- 1217 Nagy, L.G., Merenyi, Z., Hegedus, B, Balint, B. (2020). **Novel phylogenetic methods are**
1218 **needed for understanding gene function in the era of mega-scale genome**
1219 **sequencing.** *Nucl. Acids. Res.* **48**, 2209-2219. <https://doi.org/10.1093/nar/gkz1241>
1220
- 1221 Babu, M. *et al.* (2014). **Quantitative genome-wide genetic interaction screens reveal**
1222 **global epistatic relationships of protein complexes in *Escherichia coli*.** *PLoS Genetics*
1223 **10**, e1004120. <https://doi.org/10.1371/journal.pgen.1004120>
1224
- 1225 Babu, M. *et al.* (2018). **Global landscape of cell envelope protein complexes in**
1226 ***Escherichia coli*.** *Nat. Biotechnol.* **36**, 103-112. <https://doi.org/10.1038/nbt.4024>
1227
- 1228 Hu, P. *et al.* (2009). **Global functional atlas of *Escherichia coli* encompassing**
1229 **previously uncharacterized proteins.** *PLoS. Biol.* **4**, e100096.
1230 <https://doi.org/10.1371/journal.pbio.1000096>
1231
- 1232 Franceschini, A., Von Mering, C., Jensen, L.J. (2016). **SVD-phy: improved prediction of**
1233 **protein functional associations through singular value decomposition of phylogenetic**
1234 **profiles.** *Bioinformatics* **32**, 1085-1087. <https://doi.org/10.1093/bioinformatics/btv696>
1235

- 1236 Mistry J, Finn R, (2007) **Pfam: a domain-centric method for analyzing proteins and**
1237 **proteomes.** *Methods Mol. Biol.* **396**:43-58. https://doi.org/10.1007/978-1-59745-515-2_4
1238
- 1239 Huang, D.W., Sherma, B.T., Lempicki, R.A. (2009). **Bioinformatics enrichment tools:**
1240 **paths towards the comprehensive functional analysis of large gene lists.** *Nucleic Acids*
1241 *Res.* **37**, 1-13. <https://doi.org/10.1093/nar/gkn923>
1242
- 1243 Huang, D.W., Sherma, B.T., Lempicki, R.A. (2009). **Systematic and integrative analysis**
1244 **of large gene lists using DAVID Bioinformatics Resources.** *Nature Protoc.* **4**, 44-57.
1245 <https://doi.org/10.1038/nprot.2008.211>
1246
- 1247 Paul, K., Nieto, V., Carlquist, W.C., Blair, D.F., Harshey, R. (2010). **The c-di-GMP binding**
1248 **protein YcgR controls flagellar motor direction and speed to affect chemotaxis by a**
1249 **“Backstop Brake” mechanism.** *Mol. Cell.* **38**, 128-139.
1250 <https://doi.org/10.1016/j.molcel.2010.03.001>
1251
- 1252 Walker, S.L., Redman, J.A., Elimelech, M. (2004). **Role of cell surface**
1253 **lipopolysaccharides in *Escherichia coli* K12 adhesion and transport.** *Langmuir* **18**,
1254 7736-7746 (2004). <https://doi.org/10.1021/la049511f>
1255
- 1256 Alm R, Mattick J (1995) **Identification of a gene, pilV, required for type 4 fimbrial**
1257 **biogenesis in *Pseudomonas aeruginosa*, whose product possess a pre-pilin-like**
1258 **leader sequence.** *Mol. Microbiol.* **16**:485-496. [https://doi.org/10.1111/j.1365-](https://doi.org/10.1111/j.1365-2958.1995.tb02413.x)
1259 2958.1995.tb02413.x
1260
- 1261 Little A, *et al.* (2018) ***Pseudomonas aeruginosa* AlgR phosphorylation status**
1262 **differentially regulates pyocyanin and pyoverdine production.** *mBio.* **9**:e02318-17.
1263 <https://doi.org/10.1128/mBio.02318-17>.
1264
- 1265 Kearns D, Robinson J, Shimkets L, (2001) ***Pseudomonas aeruginosa* exhibits directed**
1266 **twitching motility up phosphatidylethanolamine gradients.** *J Bacteriol.* **183**:763-7.
1267 <https://doi.org/10.1128/JB.183.2.763-767.2001>
1268
- 1269 Rashid M, Kornberg A, (2000) **Inorganic polyphosphate is needed for swimming,**
1270 **swarming, and twitching motilities of *Pseudomonas aeruginosa*.** *Proc. Natl. Acad. Sci.*
1271 **97**: 4885-4890. <https://doi.org/10.1073/pnas.060030097>
1272
- 1273 Wigner, E. P., (1967) **Random matrices in physics.** *SIAM Rev.*, **9**(1):1-23.
1274 <https://doi.org/10.1137/1009001>
1275
- 1276 Cover and Thomas (2006) **Elements of information theory, 2nd edition.** ISBN: 978-0-471-
1277 24195-9
1278
1279
1280
1281
1282
1283

1284 **Figures and Figure Supplements**

1285

1286



1287

1288

1289 **Figure 1.**

1290

1291 **The SVD spectrum of OGG variation organized covariation according to biological scale.**

1292 **(A)** The data matrix, D^{OGG} , contained the number of annotations of each of 10,177 orthologous

1293 gene groups (OGGs) within each of 7,047 UniProt bacterial reference proteomes. **(B)**

1294 Benchmarks were assembled to represent prior knowledge of phylogenetic relationships

1295 (Phylogeny), indirect PPIs (Indirect PPI), and direct PPIs (Direct PPI). For each benchmark,

1296 black circles indicate the types of information represented. **(C)** Cumulative distribution functions

1297 for the mutual information (MI cdfs) shared between the SVD components of D^{OGG} and the

1298 benchmarks in panel B. Colors reflect the scheme in color legend of panel B. Shaded regions

1299 indicate ± 2 standard deviations surrounding the mean value for bootstraps of the benchmark.

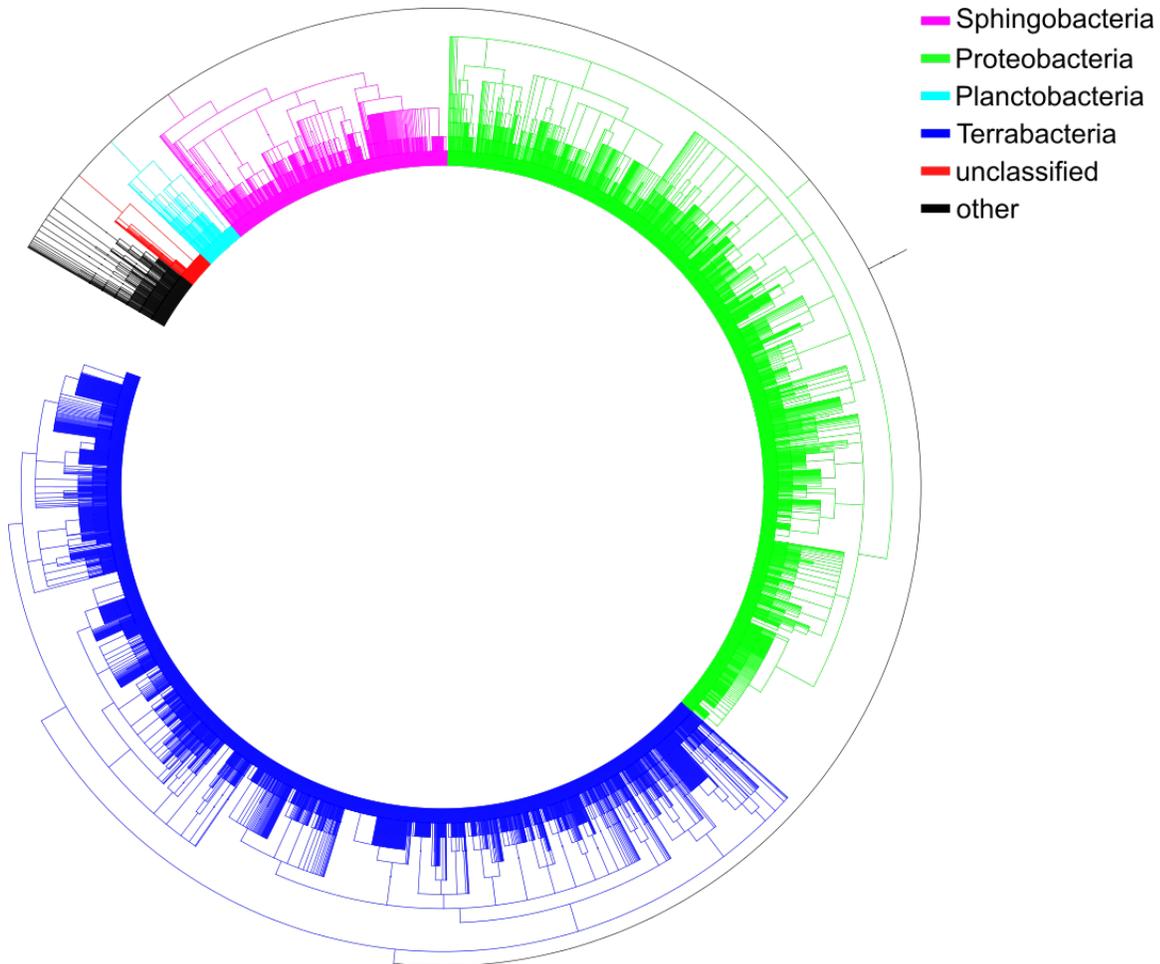
1300 Components of covariation explain progressively less of the overall data variance with

1301 increasing spectral position.

1302

1303

1304



1305

1306

1307

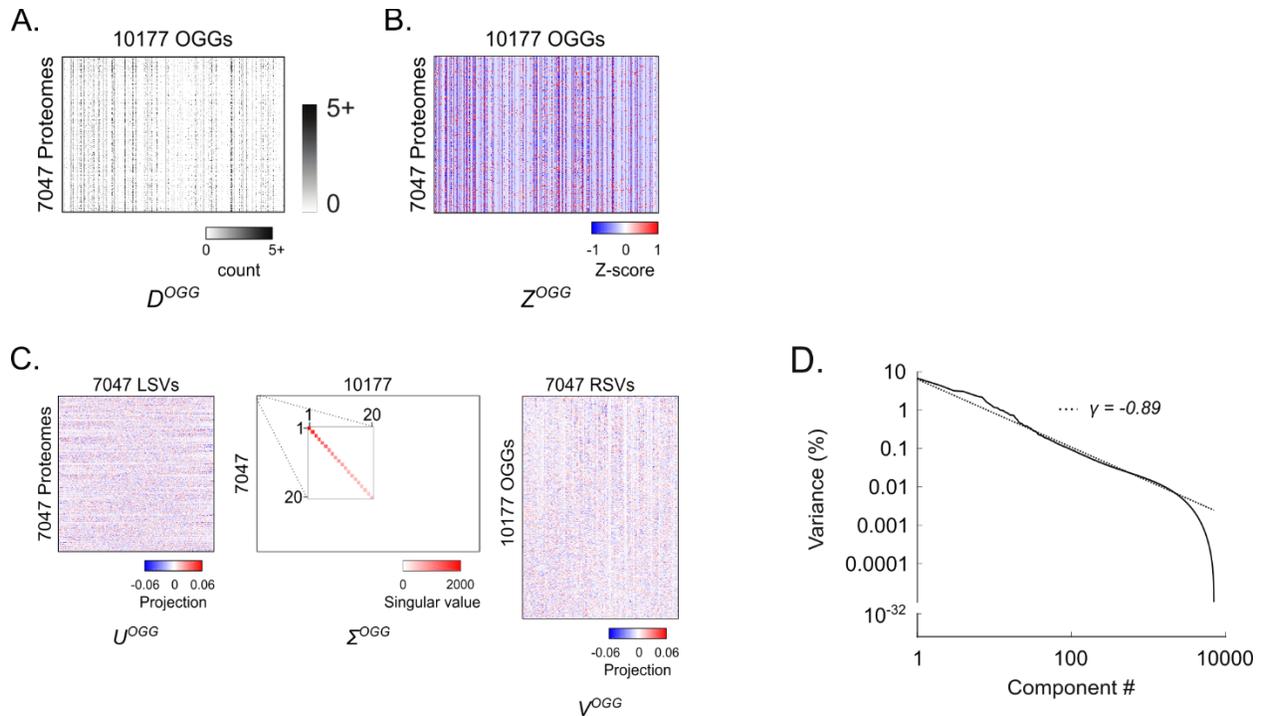
1308 **Figure 1 – figure supplement 1**

1309

1310 Phylogenetic tree of the 7,047 bacterial species represented in the UniProt reference proteomes
1311 database (release 02/2020) that served as the substrate for D^{OGG} (**Figure 1A**). The tree was
1312 generated using PhyloT based on the NCBI taxonomy database and visualized using iTOL
1313 (Letunic *et al.*, 2006).

1314

1315



1316
1317

1318 **Figure 1 – figure supplement 2**

1319

1320 **Using SVD to spectrally decompose OGG covariation in bacteria.** (A) The raw OGG
 1321 content matrix, D^{OGG} . (B) The z-scored OGG content matrix, Z^{OGG} produced by subtracting the
 1322 column mean from each element in D^{OGG} and then dividing by the column standard deviation.
 1323 (C) Application of SVD to Z^{OGG} produced three matrices: U^{OGG} , Σ^{OGG} , and V^{OGG} . U^{OGG} contains
 1324 the left singular vectors (LSVs) describing the projections of each proteome onto the SVD
 1325 components. Σ^{OGG} is a diagonal matrix with decreasing diagonal elements containing the
 1326 singular values (inset expands top 20 singular values). The magnitude of each singular value is
 1327 related to the fraction of the overall data variance explained by the corresponding SVD
 1328 component. V^{OGG} contains the right singular vectors (RSVs) describing the projections of each
 1329 OGG onto the SVD components. (D) Percent variance explained per component versus
 1330 component number. Dotted line represents the fit to a power-law distribution with the indicated
 1331 exponent (γ).
 1332

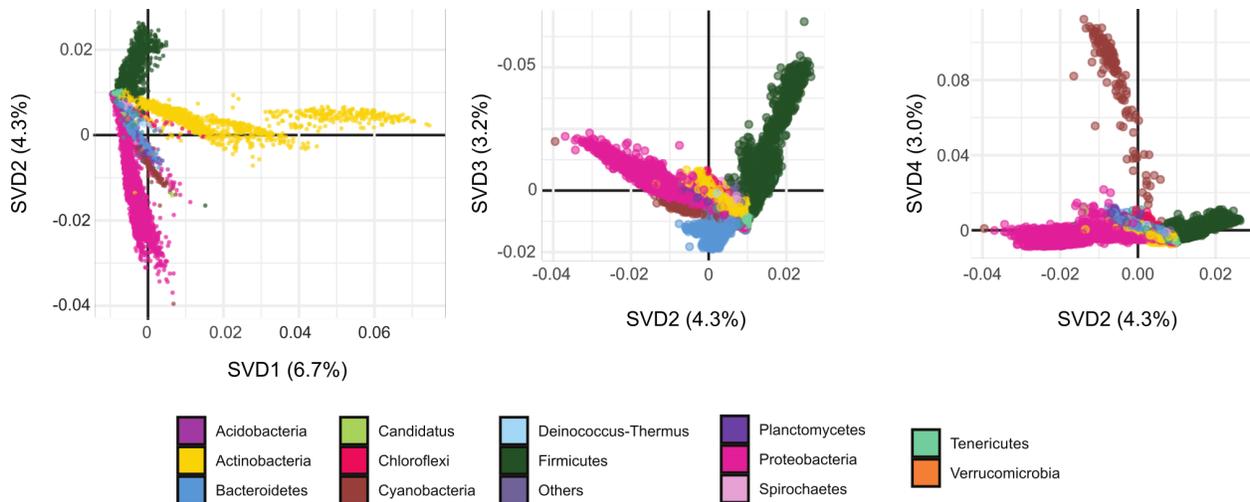


Figure 1 – figure supplement 3

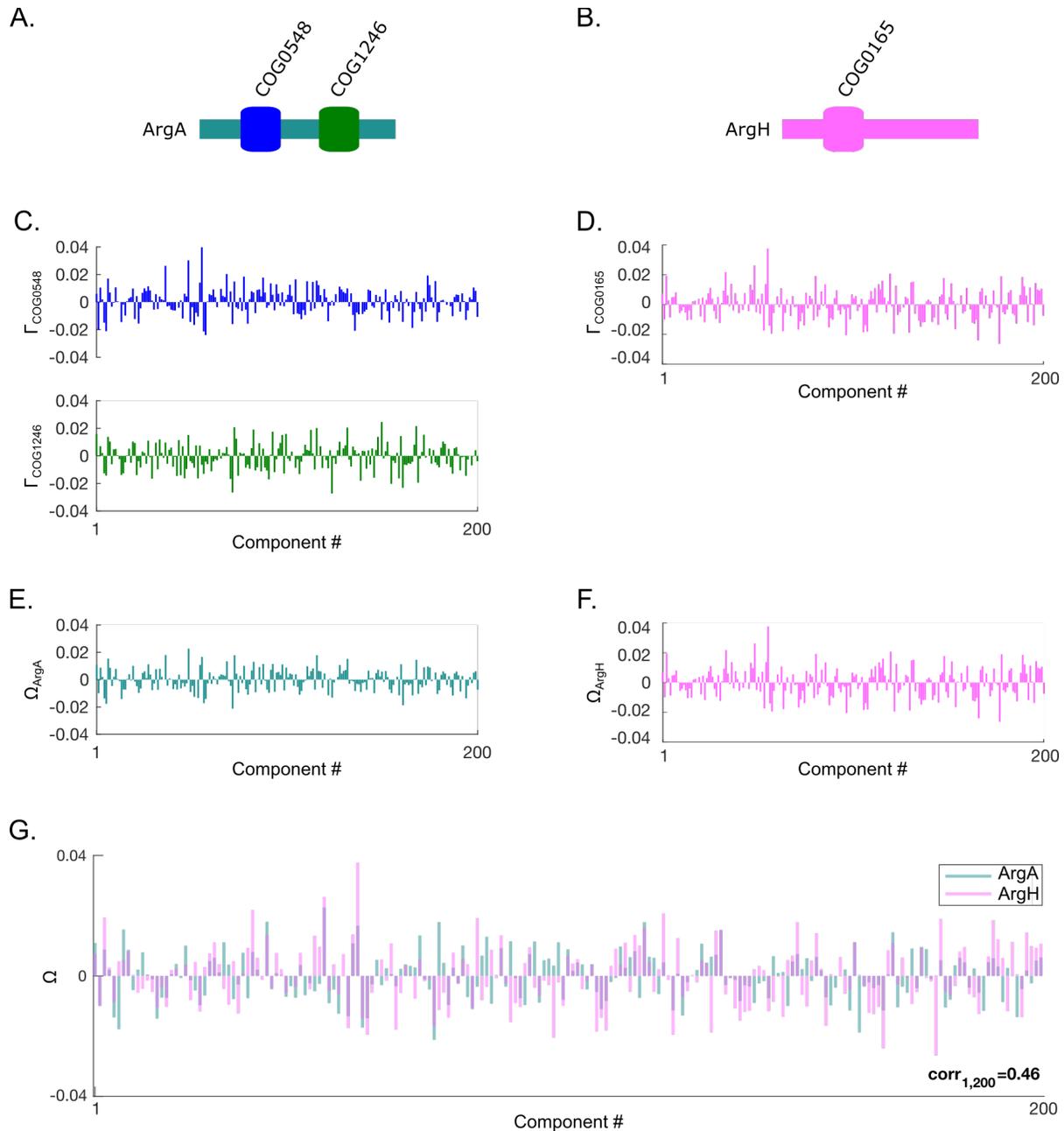
Projection of proteomes onto SVD₁, SVD₂, SVD₃, and SVD₄ of D^{OGG} colored by phylum.

1338 Proteome projections onto the SVD components are derived from U^{OGG} , the matrix of left
1339 singular vectors (LSVs) defined by applying SVD to D^{OGG} (**Figure 1 – figure supplement 2C**).

1340 Each marker represents a single proteome and is colored according by Phylum as indicated.

1341 The percent variance explained per SVD component is indicated in parentheses.

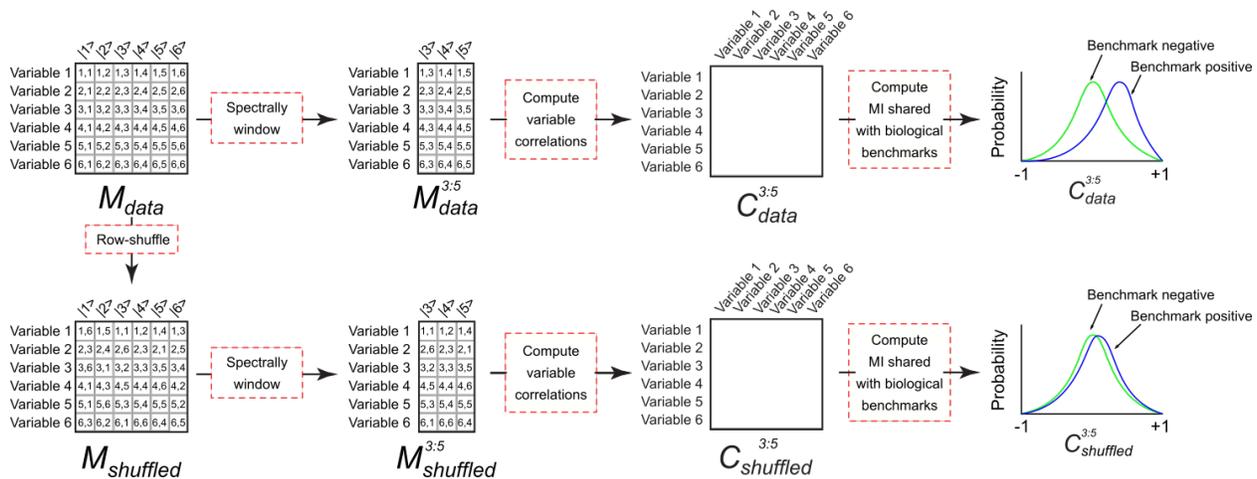
1342



1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355

Figure 1 – figure supplement 4

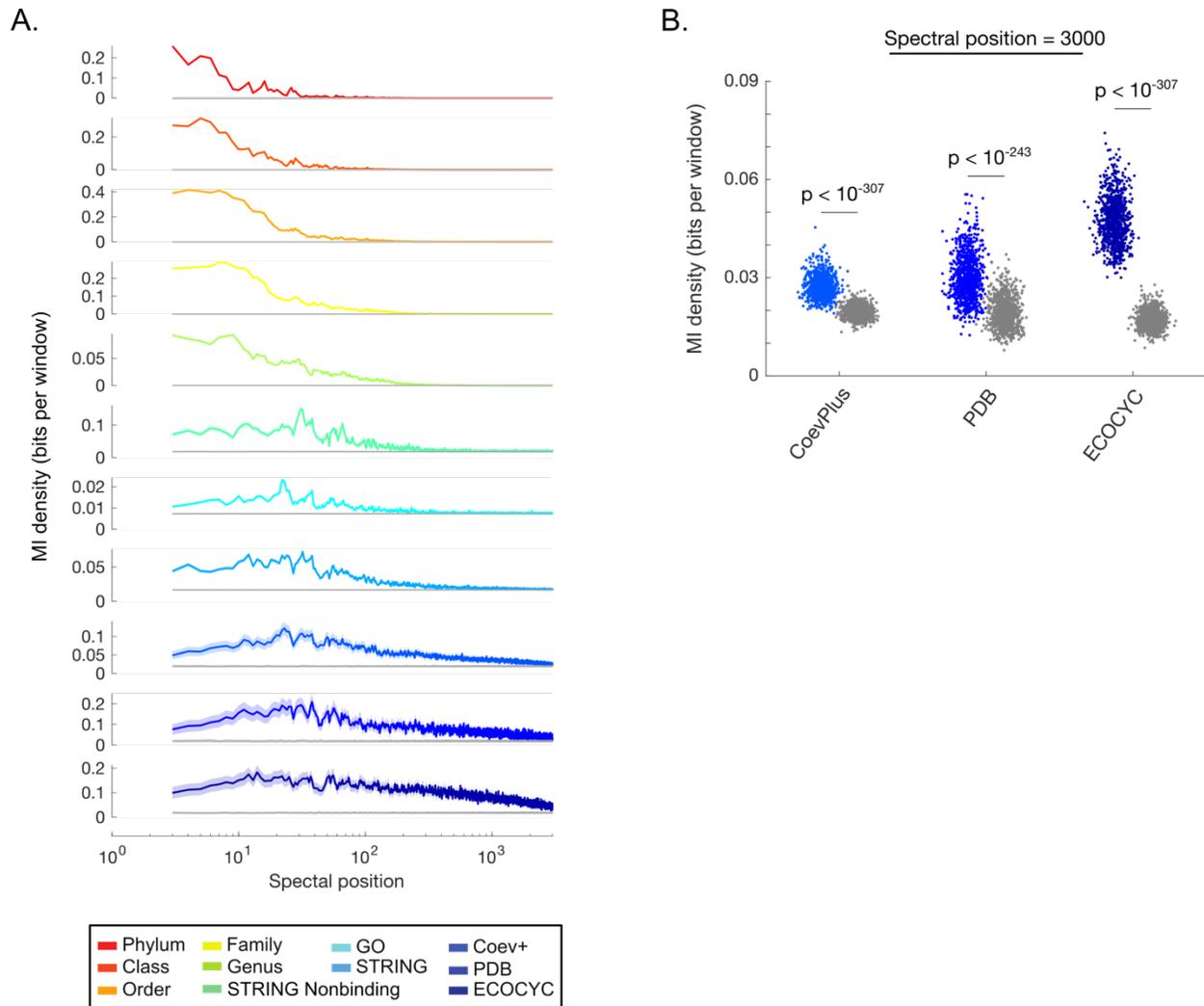
The process of computing spectral correlations between two proteins within a specific region of the SVD spectrum of D^{OGG} . (A,B) The OGG structures of *E. coli* K12 ArgA (A) and ArgH (B). (C,D) The projections (Γ) of the OGGs encoded in ArgA (COG0548 and COG1246) (C) and ArgH (COG0165) (D) onto SVD₁ to SVD₂₀₀ of the SVD spectrum of D^{OGG} . (E,F) The approximated projections (Ω) of ArgA (E) and ArgH (F) derived by averaging the projections for the OGGs encoded within each protein. (G) Overlay of the approximated protein projections of ArgA and ArgH. These two protein project similarly across SVD₁ to SVD₂₀₀ resulting in a positive spectral correlation value (Pearson correlation value shown).



1356
1357
1358
1359
1360

Figure 1 – figure supplement 5.

1361 **Estimating the MI shared between spectral correlations and a benchmark.** (Top)
 1362 Hypothetical projection matrix M_{data} consists of projections of six variables onto six SVD
 1363 components. If the variables correspond to the rows or columns of D , the initial data matrix, then
 1364 the matrix M_{data} corresponds to the complete U or V matrices produced by application of SVD to
 1365 D , respectively. M_{data} is windowed to components 3 to 5 to produce $M_{data}^{3:5}$ by discarding all
 1366 columns outside of this range. Next Pearson correlations are computed between all pairs of
 1367 rows in $M_{data}^{3:5}$ to produce the windowed spectral correlation matrix $C_{data}^{3:5}$. The MI shared
 1368 between $C_{data}^{3:5}$ and a benchmark reflects the degree to which the distribution of spectral
 1369 correlation values in $C_{data}^{3:5}$ differs for variable pairs that interact within the benchmark
 1370 (benchmark positive) and variable pairs that do not interact in the benchmark (benchmark
 1371 negative). (Bottom) To estimate the MI produced by spurious correlations, M_{data} is subjected to
 1372 random row permutation to generate $M_{shuffled}$. This process maintains the distribution of
 1373 projection values for each variable but erases biologically meaningful spectral correlations
 1374 leaving only the spurious correlations. $M_{shuffled}$ is subjected to the identical windowing, row
 1375 correlation computation, and MI calculations as described above for M_{data} . The biologically
 1376 meaningful MI is estimated to be the difference between the MI estimate for $C_{data}^{3:5}$ and
 1377 $C_{shuffled}^{3:5}$.
 1378

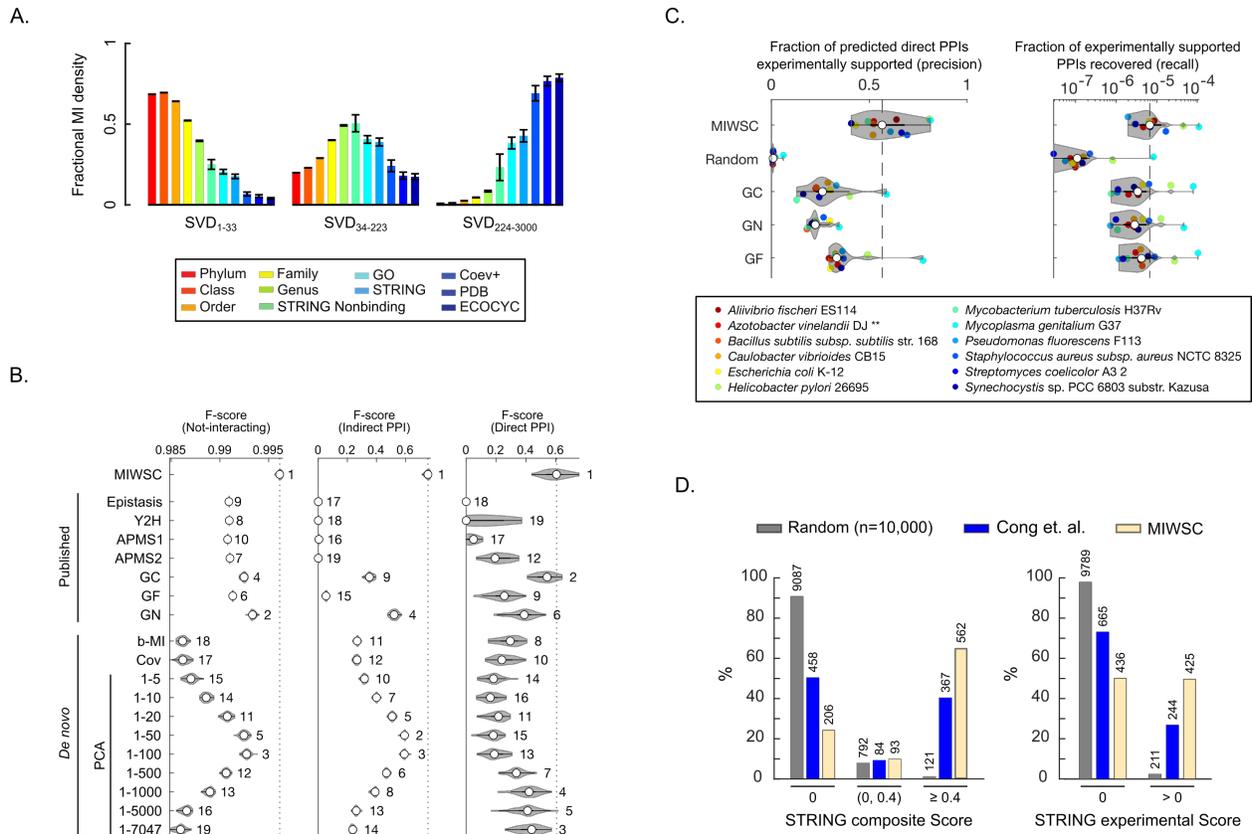


1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395

Figure 1 – figure supplement 6

Quantifying the biological information contained within different regions of the SVD spectrum of D^{OGG} . (A) MI density contained within a 5-component window versus spectral position for all windows between SVD_1 and SVD_{3000} in the SVD spectrum of D^{OGG} . Colored and gray lines represent the MI values for the data matrix or the model of spurious spectral correlations, respectively. Lines and shaded contours represent the mean ± 2 standard deviations for the bootstraps of each benchmark. (B) MI density contained within SVD_{2995} - SVD_{3000} of the SVD spectrum of D^{OGG} . Each dot represents the MI value for a single bootstrap of the indicated benchmark. Colored dots represent the MI for the data matrix and gray dots represent the MI for the model of spurious correlations. P-values produced by a paired Student's t-test are indicated.

1396



1397

1398

1399

1400

1401

1402

1403

1404

1405

1406

1407

1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

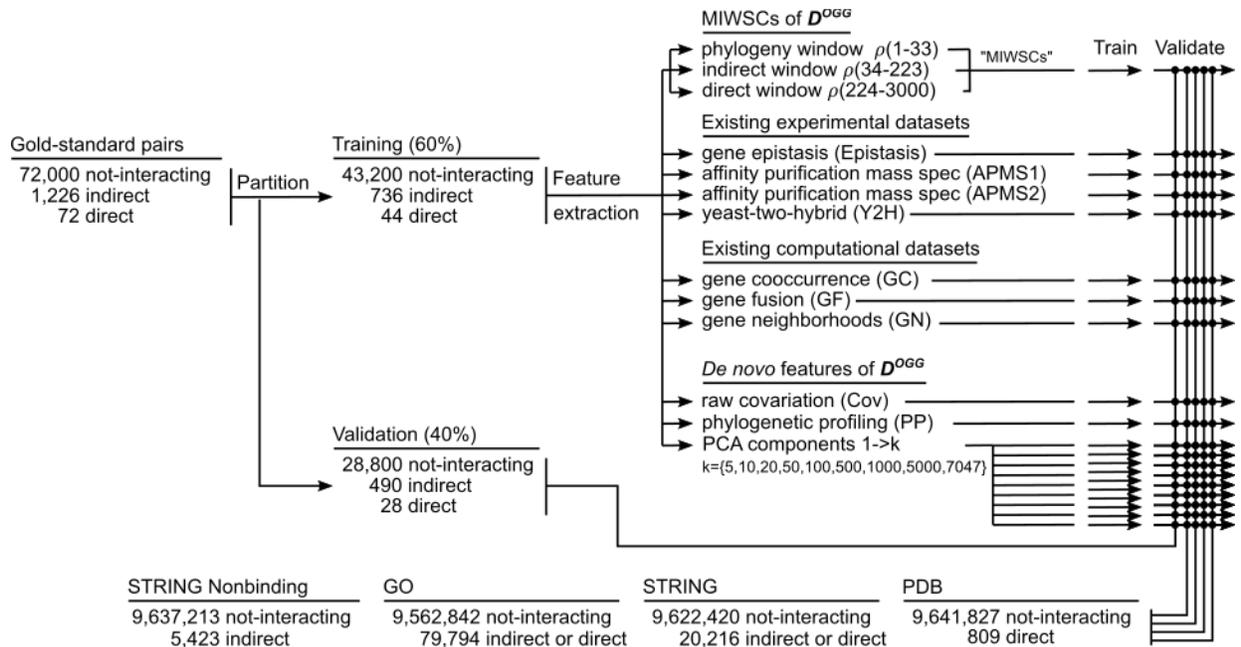
Figure 2. MI windowed spectral correlations (MIWSCs) enable accurate classification of protein pairs as either not-interacting, indirect PPI, or direct PPI. (A) Fractional MI density regarding each benchmark contained within spectral correlations computed across SVD₁₋₃₃, SVD₃₄₋₂₂₃, or SVD₂₃₄₋₃₀₀₀ of D^{OGG} . Color scheme is defined in the legend and follows that of **Figure 1B,C.** **(B)** F-scores for predicting interaction classes for pairs in an independent validation set of *E. coli* K12 proteins using RF models trained on either MIWSCs, quantitative features of published datasets derived from experimental methods (gene epistasis [epistasis], yeast-two-hybrid [Y2H], affinity purification mass spectrometry [APMS1, APMS2]), quantitative features of published datasets derived from computational methods (gene cooccurrence [GC], gene fusion [GF], gene neighborhood [GN]), or quantitative features of established computational methods derived *de novo* from D^{OGG} (binary mutual information [b-MI], covariation [Cov], Principle Components Analysis including components 1-k [PCA]). The violin plots describe the distribution of F-scores for models trained and validated on 50 random partitions of the gold-standard dataset (**Figure 2 – figure supplement 1**). Numbering indicates the rank of the median F-score for models trained on each feature (**Table S4**). **(C)** Precision (left) and recall (right) for direct PPI predictions in 12 phylogenetically diverse organisms using RF models trained on the MIWSCs of *E. coli* K12 proteins benchmarked against the experimentally supported PPIs in the STRING database. Comparisons are made to a set of 10,000 randomly selected pairs and to the 'high confidence' predictions in the STRING database subchannels for the methods of gene cooccurrence (GC), gene neighborhood (GN), and gene fusion (GF). Vertical dashed line indicates the median value for RF models trained on MIWSCs. ** in legend indicates an organism that was not part of the input dataset D^{OGG} (**Table S4**). **(D)** Percent of predicted direct PPIs in *M. tuberculosis* H37Rv supported by an absent (0),

1423 low (0 to 0.4), or high (>0.4) composite score (left) or an absent (0) or present (>0) experimental
1424 subchannel score (right) in the STRING database. Comparisons were made between the
1425 methods of random selection (Random), amino acid coevolution ('Cong *et al.*', Cong *et al.*,
1426 2019), or RF models trained on MIWSC features of *E. coli* K12 proteins (MIWSC). Numbers of
1427 predicted interactions in each bin are indicated (**Table S4**).

1428

1429

1430



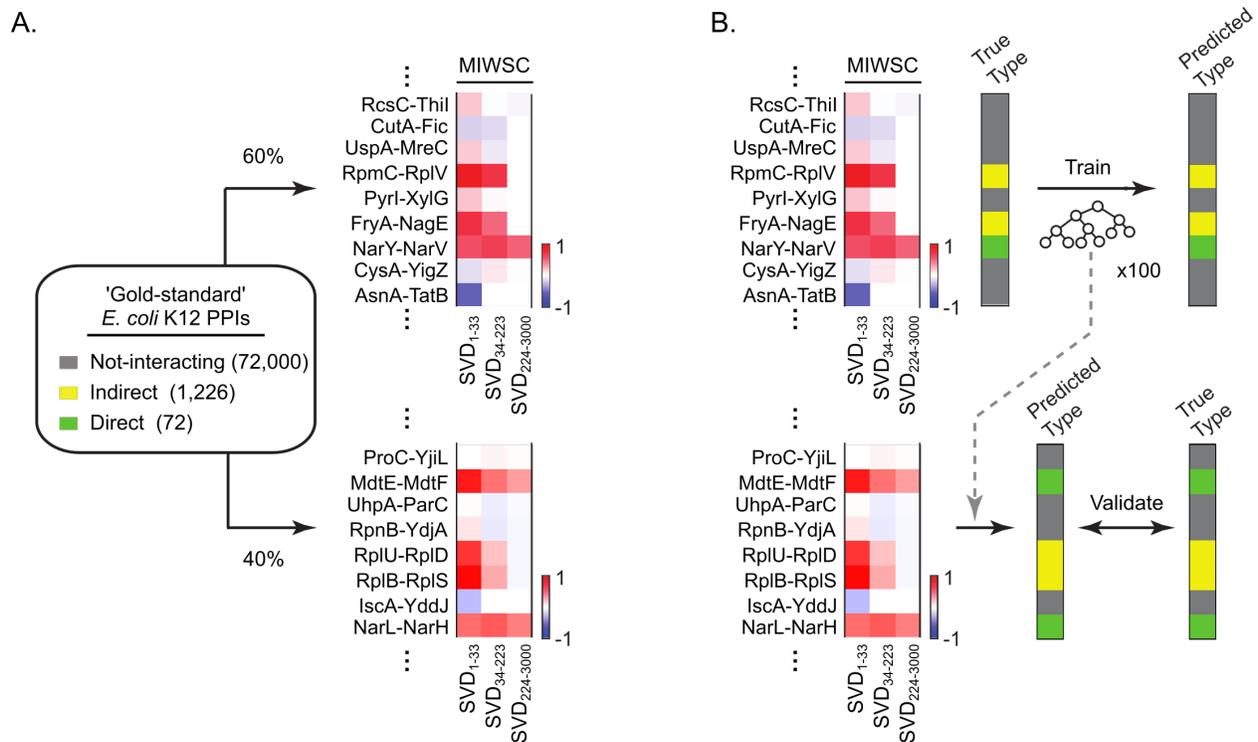
1431
1432

1433 **Figure 2 – figure supplement 1**

1434

1435 **Multi-level classification task where RF models were trained and validated for predicting**
 1436 **indirect and direct PPIs in *E. coli* K12 using spectral correlations features, existing**
 1437 **computational features, or existing experimental features.** A gold-standard dataset of well-
 1438 characterized *E. coli* K12 protein pairs was assembled and partitioned into training and
 1439 validation datasets. The labeled examples from the training set were used to train RF models to
 1440 classify protein interactions using different features as indicated. The performance of the various
 1441 RF models was benchmarked and compared by computing F-scores for classifying PPIs in the
 1442 validation dataset, out-of-bag examples in the training dataset, and PPIs in four additional
 1443 comprehensive benchmarks. This process of partitioning the gold-standard dataset, training,
 1444 and validation was repeated 50 times to evaluate the reproducibility of RF model performance.

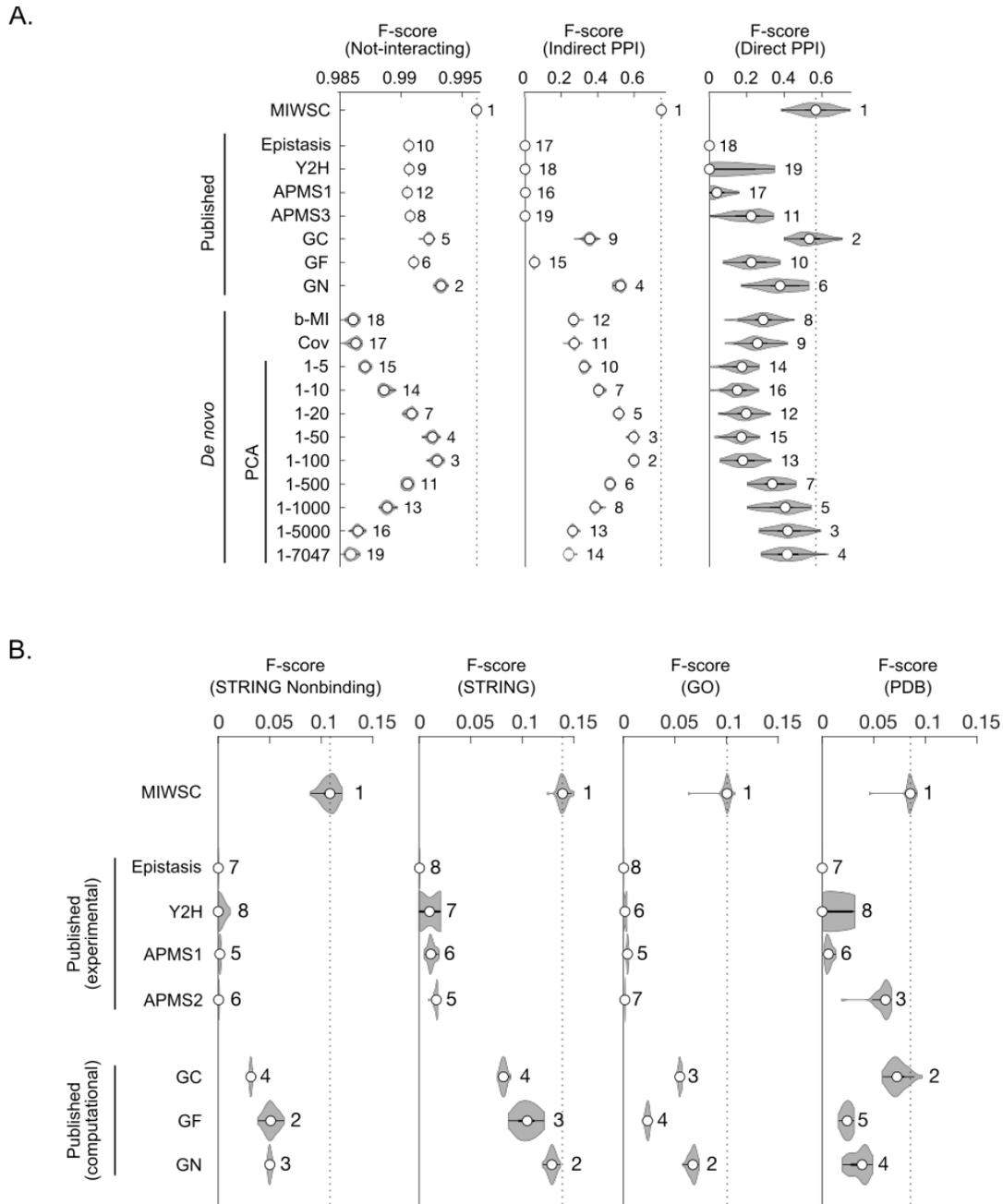
1445



1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461

Figure 2 – figure supplement 2

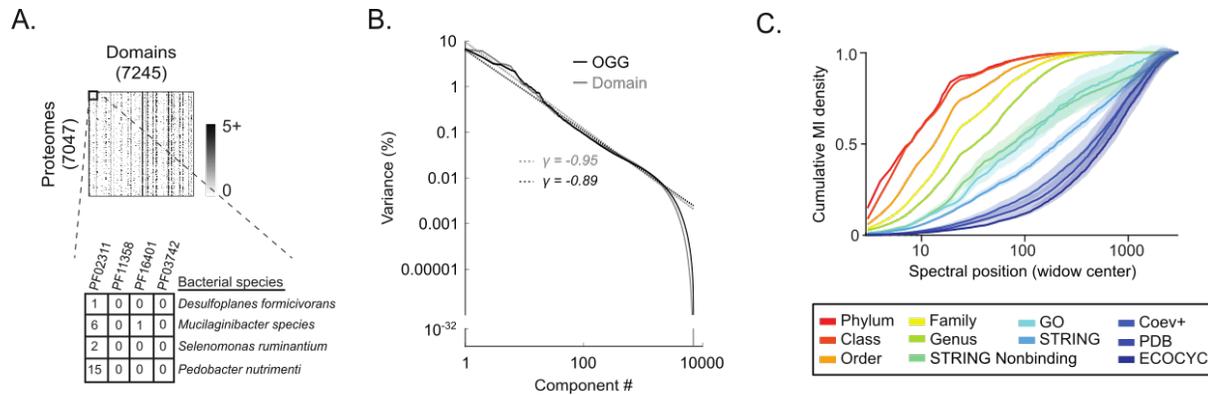
Extracting MIWSC features for the training and validation of RF models to classify *E. coli* K12 protein pairs as not-interacting ('not-interacting'), indirect PPI ('indirect'), or direct PPI ('direct'). (A) A gold standard dataset of well characterized *E. coli* K12 protein pairs was assembled and randomly partitioned into training (60%) and validation (40%) datasets. An MIWSC feature was extracted for each protein pair in the training and validation partitions of the gold standard dataset. The MIWSC feature consists of a set of three spectral correlations computed across SVD₁₋₃₃, SVD₃₄₋₂₂₃, and SVD₂₂₄₋₃₀₀₀. Each pixel of the heat map is the spectral correlation for the protein pair indicated in the row across the SVD components indicated in the column. (B) Using only the MIWSC features, RF models were trained using the labeled examples in the training dataset and then challenged to predict the unlabeled examples in the validation dataset.



1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474

Figure 2 – figure supplement 3

F-scores for predicting interaction classes for out-of-bag examples in the training datasets (A) and four additional comprehensive benchmarks (B). The violin plots describe the distribution of F-scores for models trained and validated on 50 random partitions of the gold-standard dataset (**Figure 2 – figure supplement 1**). Numbering indicates the rank of the median F-score for models trained on each feature (**Table S4**). Feature descriptions can be found in the legend of **Figure 2**.

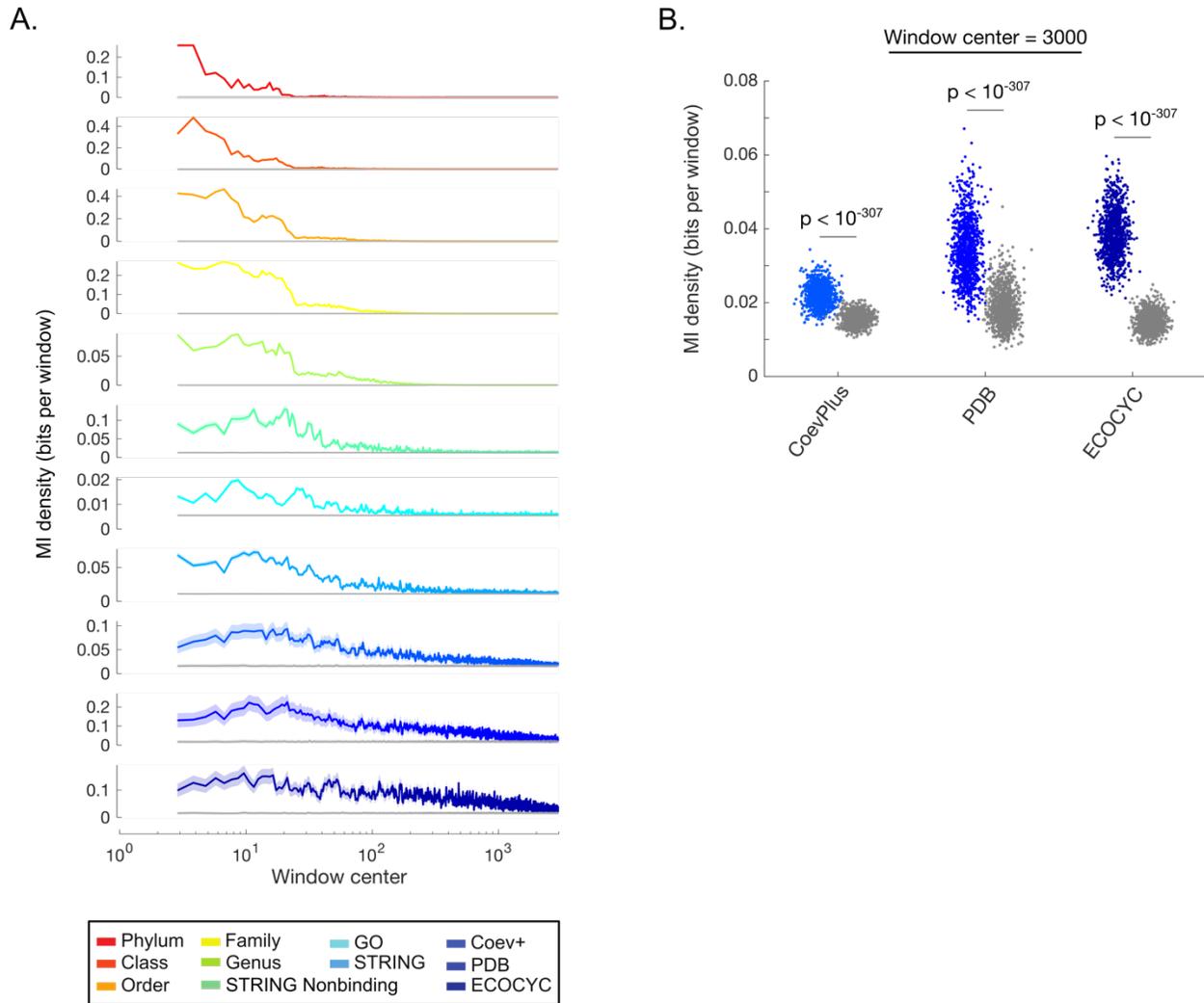


1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490

Figure 3.

The SVD spectrum of protein domain variation organized covariation according to biological scale. (A) The domain content matrix, D^{domain} , contained the number of annotations of each of 7,245 conserved protein domains within each of the 7,047 UniProt bacterial reference proteomes. **(B)** Percent variance explained per component versus component number for the SVD spectra of D^{OGG} (black) and D^{domain} (gray). Dotted lines represent a power-law distribution with the indicated exponent (γ). **(C)** Cumulative distribution functions for the mutual information (MI cdfs) shared between the SVD components of D^{domain} and the indicated benchmarks. Shaded regions indicate ± 2 standard deviations surrounding the mean value for bootstraps of each benchmark (STAR Methods).

1491



1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

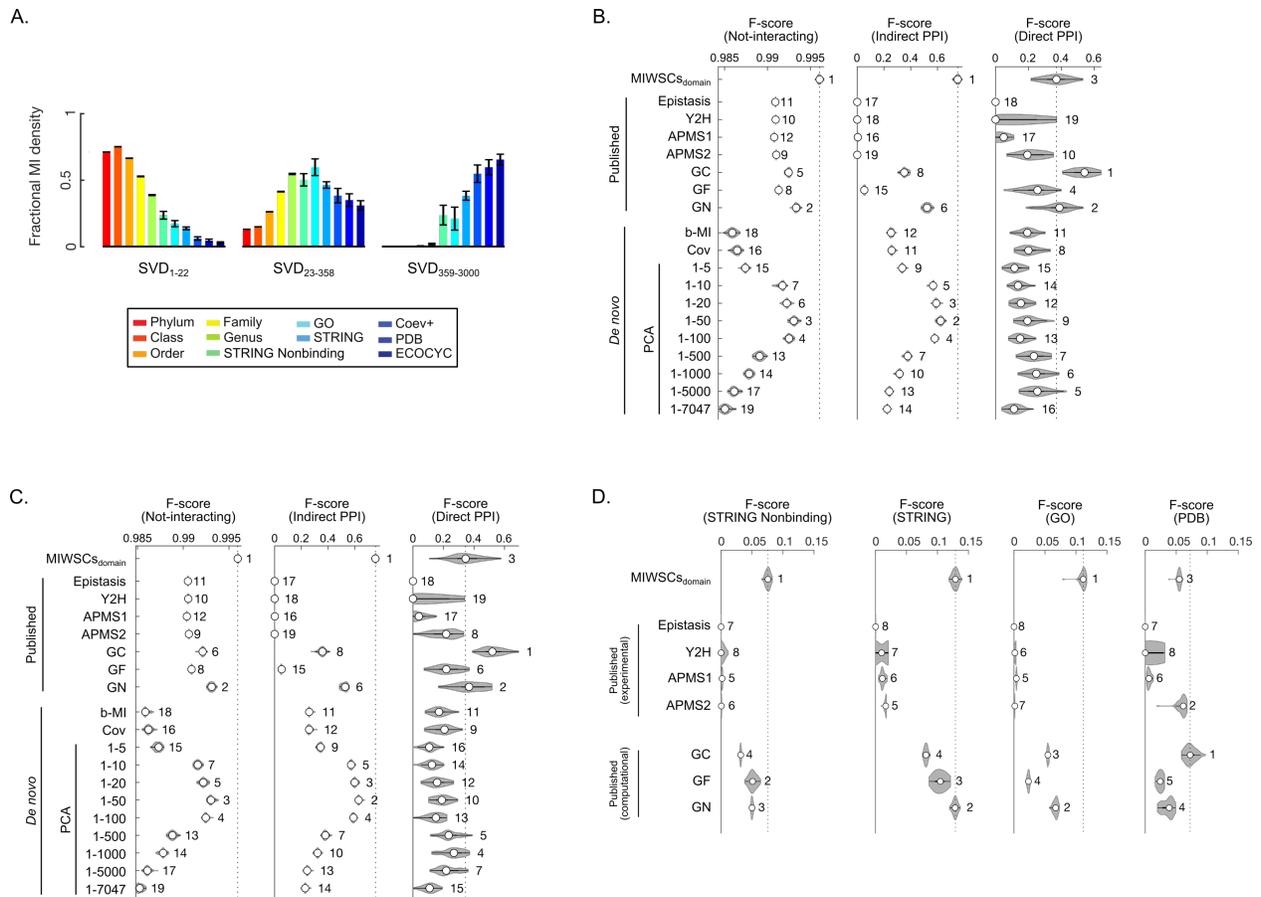
1509

1510

1511

Figure 3 – figure supplement 1

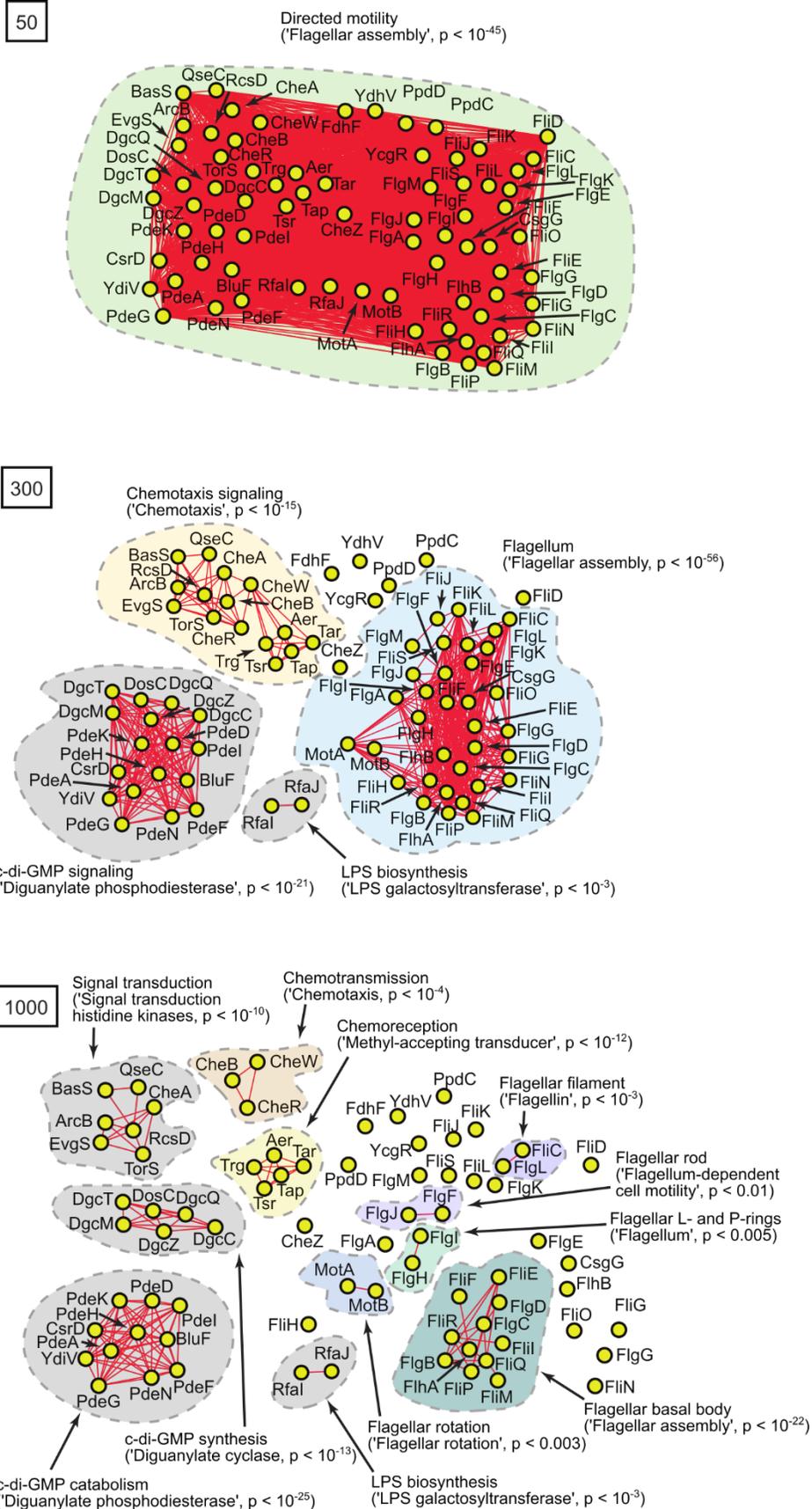
Quantifying the biological information contained within different regions of the SVD spectrum of D^{domain} . (A) MI density contained within a 5-component window versus spectral position for all windows between SVD_1 and SVD_{3000} in the SVD spectrum of D^{domain} . Colored and gray lines represent the MI values for the data matrix or the model of spurious spectral correlations, respectively. Lines and shaded contours represent the mean ± 2 standard deviations for the bootstraps of each benchmark. (B) MI density contained within SVD_{2995} - SVD_{3000} of the SVD spectrum of D^{domain} . Each dot represents the MI value for a single bootstrap of the indicated benchmark. Colored dots represent the MI for the data matrix and gray dots represent the MI for the model of spurious correlations. P-values produced by a paired Student's t-test are indicated.



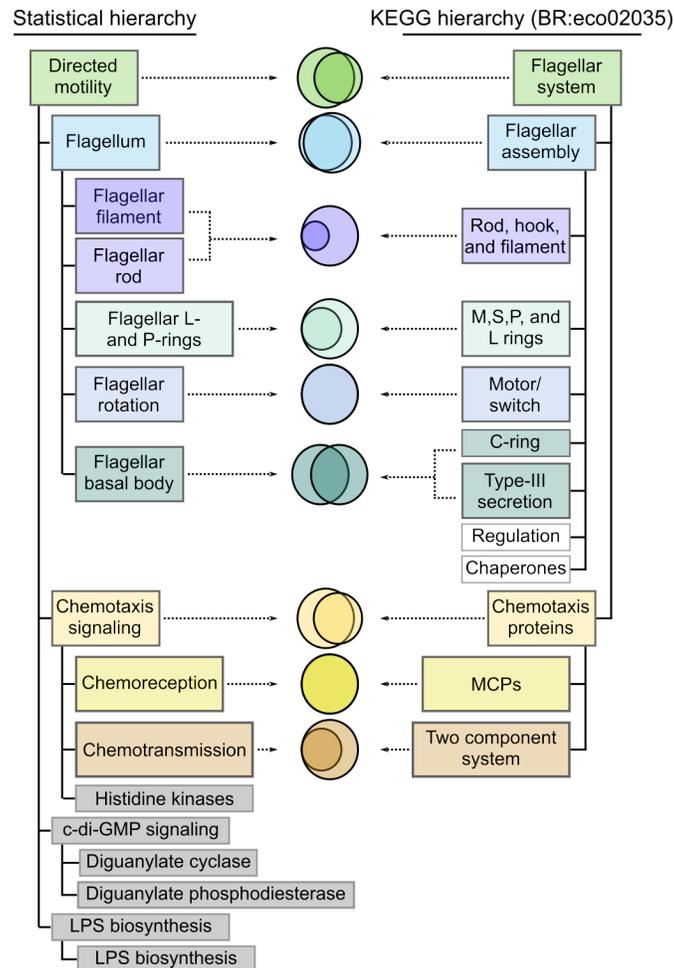
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529

Figure 3 – figure supplement 2

Domain-based MI windowed spectral correlations (MIWSCs_{domain}) enable accurate classification of protein pairs as either not-interacting, indirect PPI, or direct PPI. (A) Fractional MI density regarding each benchmark contained within spectral correlations computed across SVD₁₋₂₂, SVD₂₃₋₃₅₈, or SVD₃₅₉₋₃₀₀₀ of D^{domain} . Color scheme is defined in the legend and follows that of **Figure 1B,C**. **(B-D)** F-scores for classifying protein pairs in the validation dataset **(B)**, out-of-bag examples from the training dataset **(C)**, or additional examples in four comprehensive benchmarks **(D)**. The violin plots describe the distribution of F-scores for models trained and validated on 50 random partitions of the gold-standard dataset (**Figure 2 – figure supplement 1**). Numbering indicates the rank of the median F-score for models trained on each feature (**Table S6**). Feature descriptions can be found in the legend of **Figure 2**.



1531



1532

1533

1534

1535

1536

1537

1538

1539

1540

1541

1542

1543

1544

1545

1546

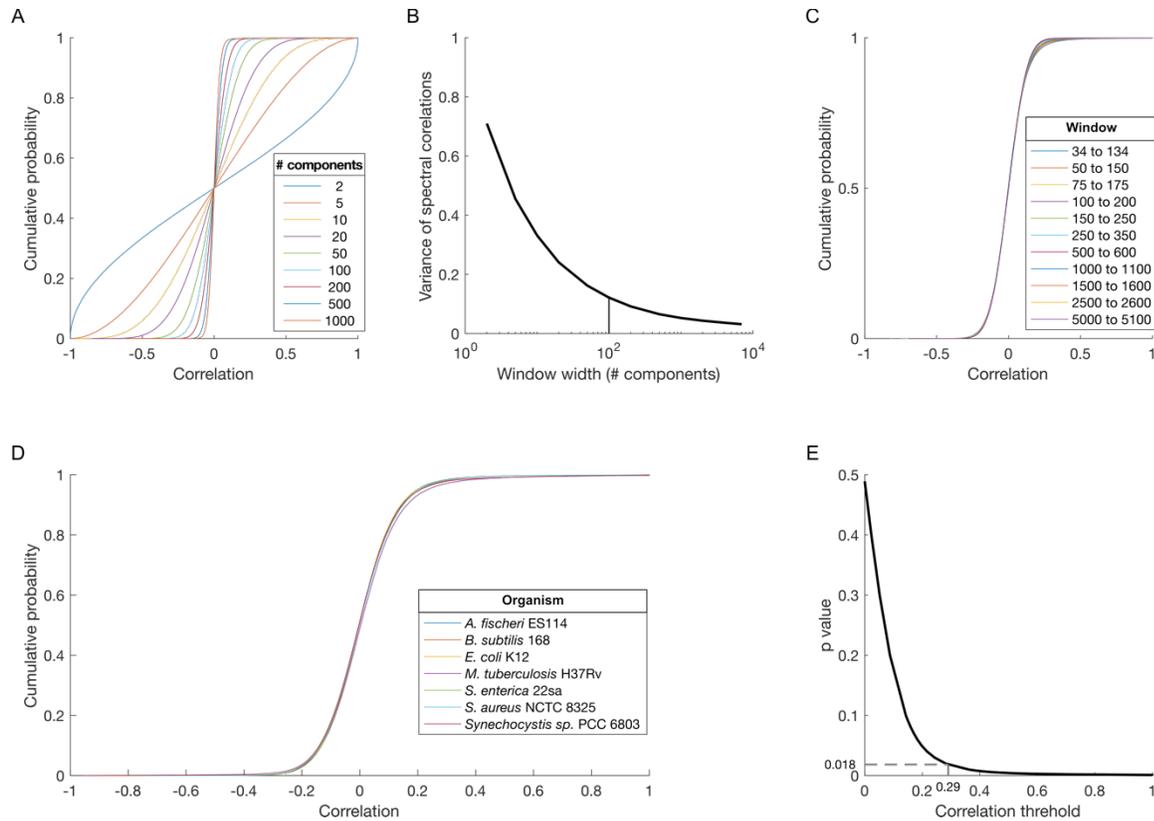
1547

1548

1549

Figure 4.

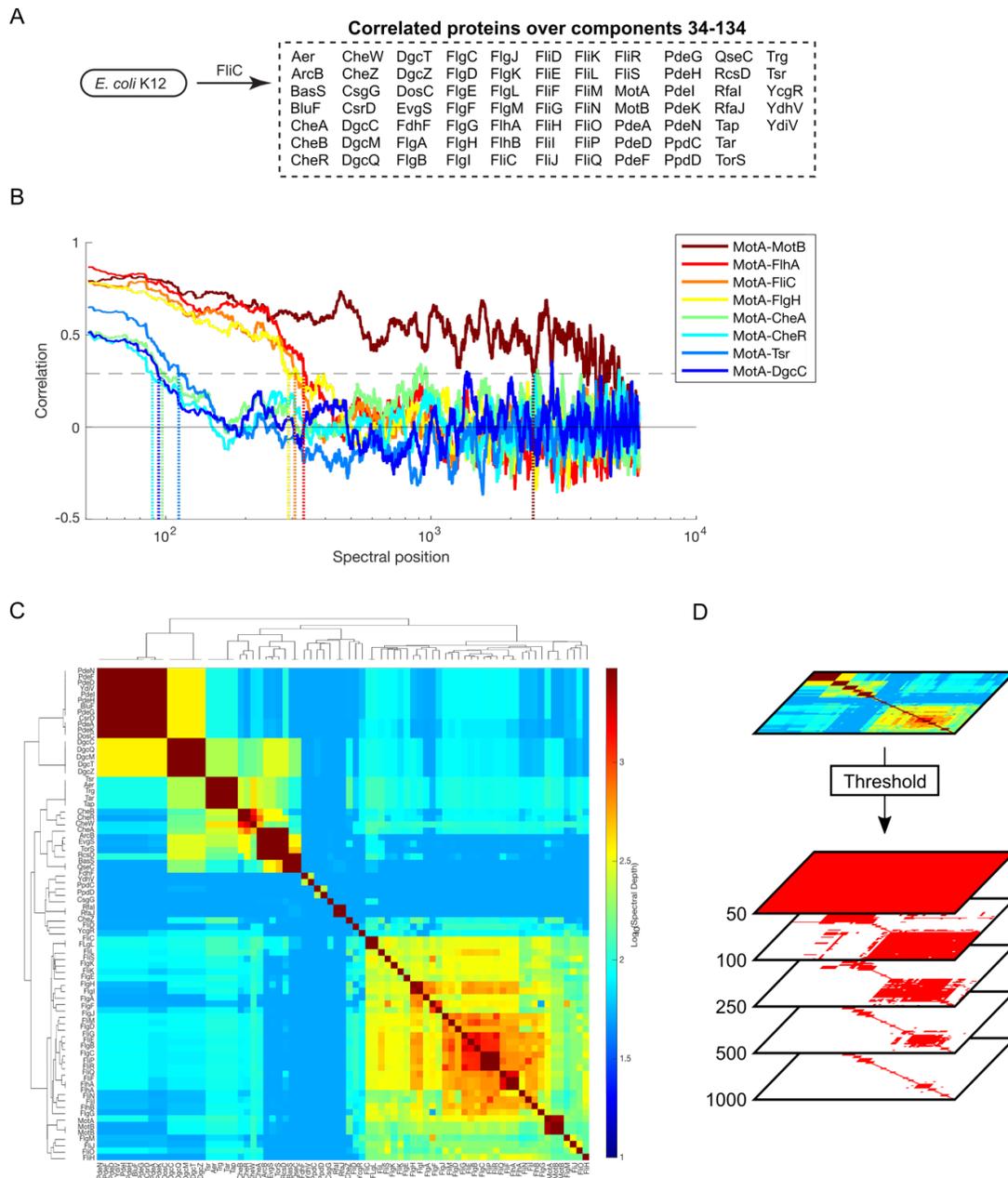
A hierarchical model of *E. coli* K12 motility derived by serially thresholding ‘spectral depth’ resembles that described in the KEGG database. (A) The model contained 75 proteins that were identified as significantly correlated with FlhC across SVD₃₄ to SVD₁₃₄. Statistical interaction networks were defined by thresholding spectral depth at 50 (top panel), 300 (middle panel), and 1000 (bottom panel). Nodes (yellow circles) represent proteins and edges (red lines) represent statistical interactions. Shaded contours identify discrete subnetworks and are labeled with their assigned function based on interpretation of gene-set enrichment analysis (GSEA) and literature review. The most significantly enriched ontological term produced by GSEA and the associated p-value is shown in parentheses for each subnetwork (**Table S7**). (B) Comparison of the statistically-derived model (left) to the KEGG model (BR:eco02035, right) of *E. coli* K12 motility. Venn diagrams represent the overlap between the sets of proteins in the indicated subnetwork of panel A (left) and the indicated KEGG category (right).



1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564

Figure 4 – figure supplement 1

Developing a null model of random protein-protein spectral correlations within the SVD spectrum of D^{OGG} . (A) Cumulative distribution functions (cdfs) for spectral correlations between all proteins in *E. coli* K12 across windows of different widths (legend) centered on component 1001. (B) Variance of the distributions in panel A plotted versus window width. Vertical line indicates a window width of 100 components. (C) cdfs for spectral correlations between all proteins in *E. coli* K12 across the indicated 100-component spectral windows (legend). (D) cdfs for spectral correlations for all proteins in proteomes from diverse organisms (legend) across the 100-component window centered on component 84. (E) p-value versus correlation threshold for spectral correlations between proteins in *E. coli* K12 across SVD₃₄ to SVD₁₃₄. The chosen correlation threshold (0.29) and associated p-value (0.018) are indicated.

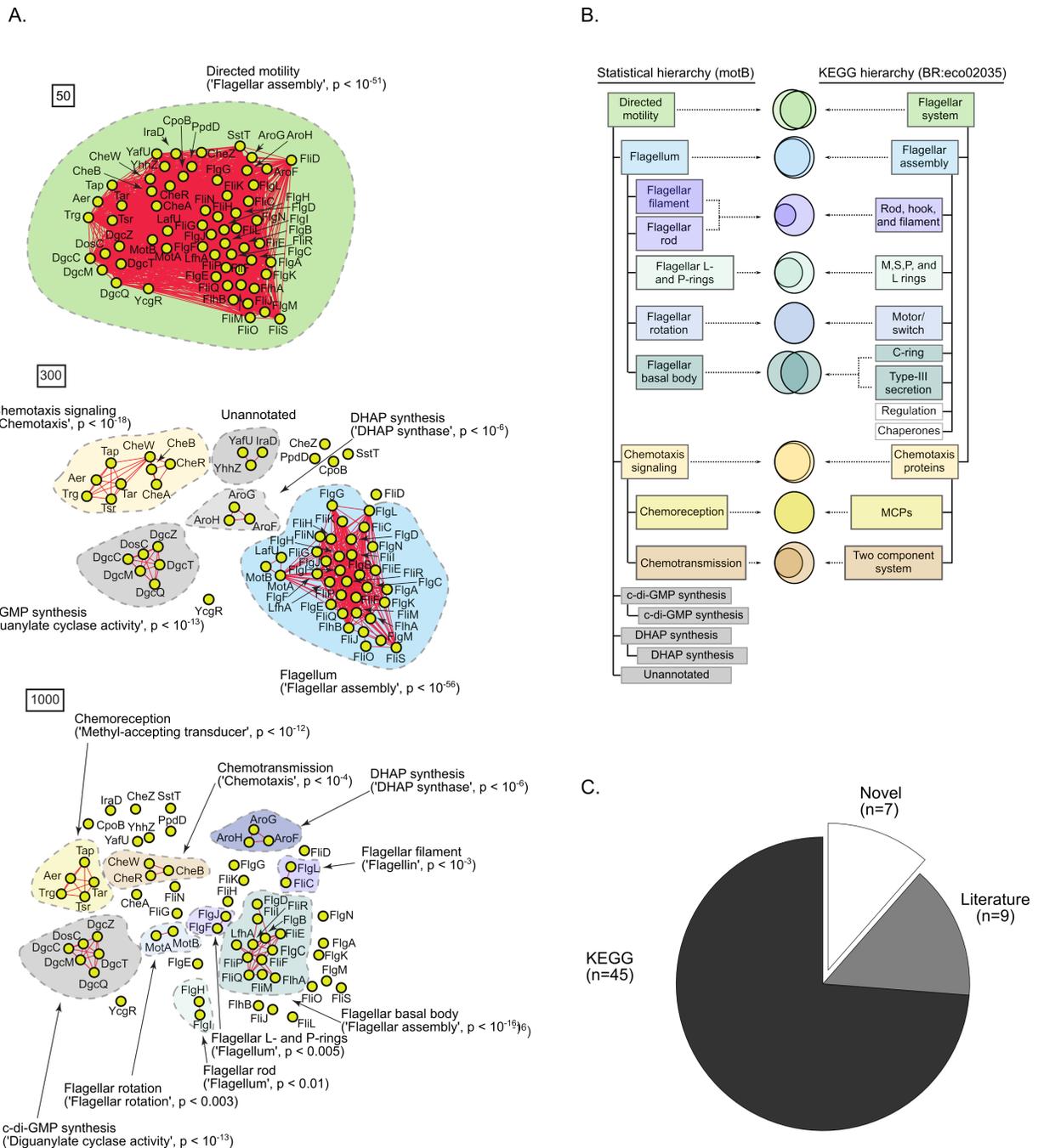


1565
1566

1567 **Figure 4 – figure supplement 2**

1568

1569 **Process for constructing the FliC interaction networks appearing in Figure 4 using**
 1570 **thresholded spectral depth. (A)** Proteins that shared significant spectral correlations with FliC
 1571 across SVD₃₄ to SVD₁₃₄. **(B)** Pairwise spectral correlations between indicated protein pairs
 1572 (legend) as a function of spectral position. Horizontal dashed line represents the threshold of
 1573 significant spectral correlation described in **Figure 4 – figure supplement 1E**. Dashed vertical
 1574 lines represent the ‘spectral depth’—the spectral position at which spectral correlation first falls
 1575 below the significance threshold. **(C)** Hierarchically clustered spectral depth matrix for all pairs
 1576 of proteins in panel A. **(D)** Set of adjacency matrices derived from thresholding the spectral
 1577 depth matrix in panel C. Red and white pixels indicate proteins that do or do not share a
 1578 spectral depth exceeding the indicated thresholds respectively.



1579
1580

1581 **Figure 4 – figure supplement 3**

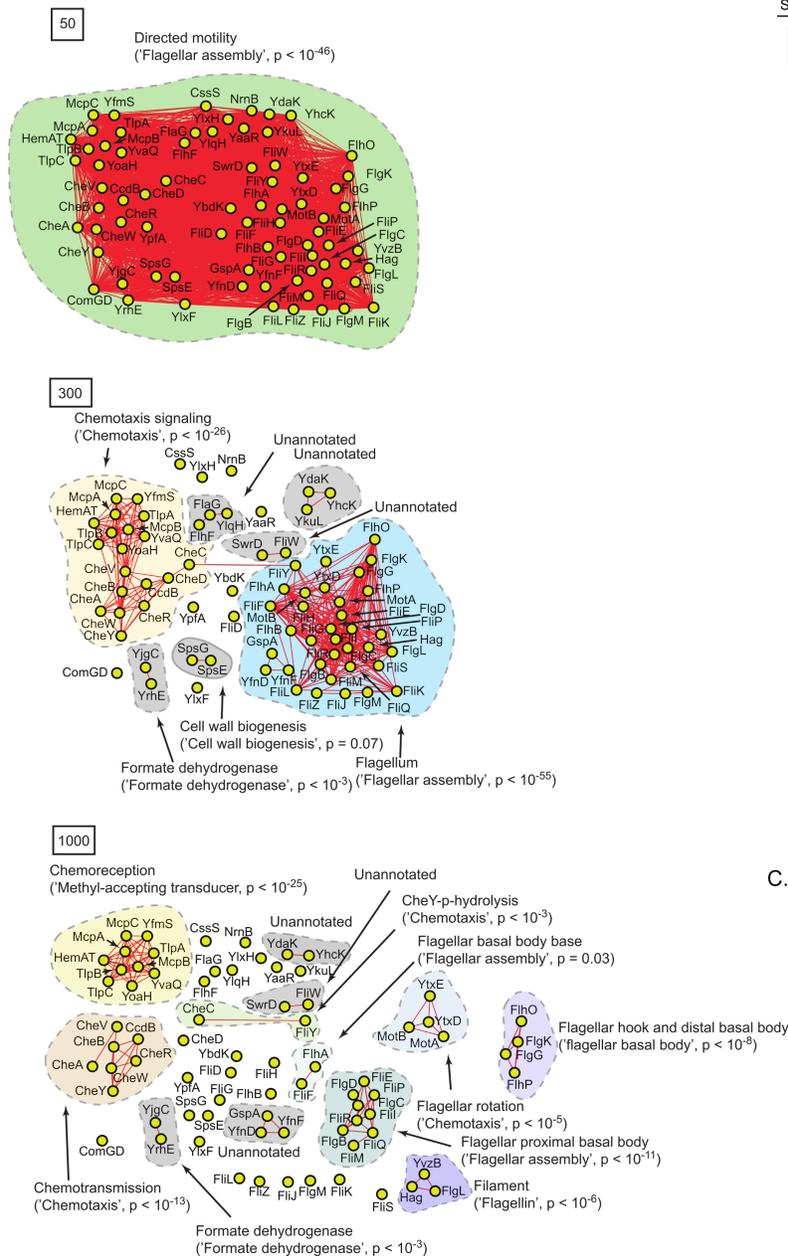
1582

1583 **A statistically-derived hierarchical model of motility in *E. coli* K12 using MotB as a query**
 1584 **protein. (A)** 60 proteins in *E. coli* K12 shared significant spectral correlations with MotB across
 1585 SVD₃₄ to SVD₁₃₄. Statistical interaction networks were defined by thresholding spectral depth at
 1586 50 (top panel), 300 (middle panel), and 1000 (bottom panel). Nodes (yellow circles) represent
 1587 proteins and edges (red lines) represent statistical interactions. Shaded contours identify
 1588 discrete subnetworks and are labeled with their assigned function based on interpretation of
 1589 gene-set enrichment analysis (GSEA) and literature review. The most significantly enriched
 1590 ontological term produced by GSEA and the associated p-value is shown in parentheses for

1591 each subnetwork (**Table S8**). B) Comparison of the statistically-derived model using MotB (left)
1592 to the KEGG model (BR:eco02035, right) of *E. coli* K12 motility. Venn diagrams represent the
1593 overlap between the sets of proteins in the indicated subnetwork of panel A (left) and the
1594 indicated KEGG category (right). (C) Pie graph of the number of proteins in the statistical model
1595 that are represented in the KEGG hierarchy ('KEGG'), missing from the KEGG hierarchy but
1596 supported by experimental evidence in the literature ('Literature'), or absent from the KEGG
1597 hierarchy and the literature ('Novel').
1598
1599

1600
1601

A.



1602
1603
1604
1605
1606
1607
1608
1609
1610
1611

Figure 4 – figure supplement 4

A statistically-derived hierarchical model of motility in *B. subtilis* 168 using Hag as a query protein. (A) 74 proteins in *B. subtilis* 168 shared significant spectral correlations with Hag over SVD₃₄ to SVD₁₃₄. Statistical interaction networks were defined by thresholding spectral depth at 50 (top panel), 300 (middle panel), and 1000 (bottom panel). Nodes (yellow circles) represent proteins and edges (red lines) represent statistical interactions. Shaded contours

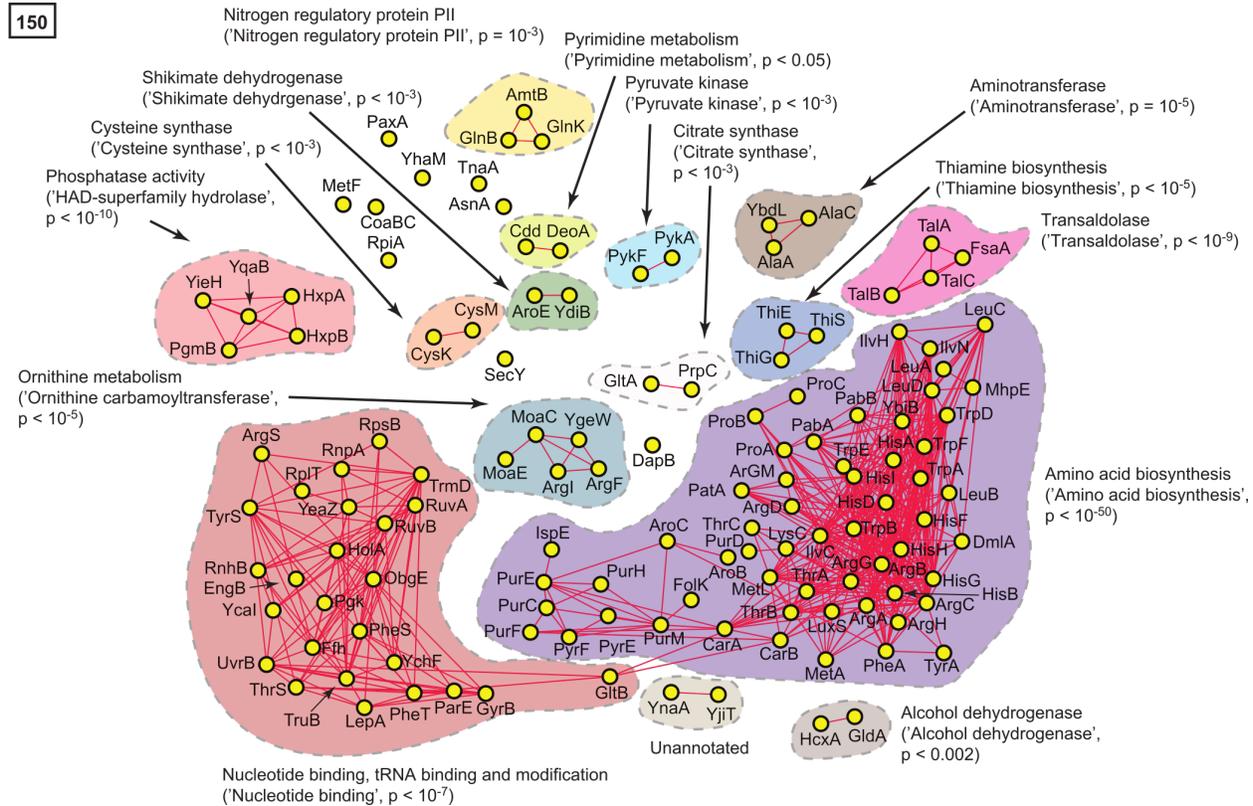
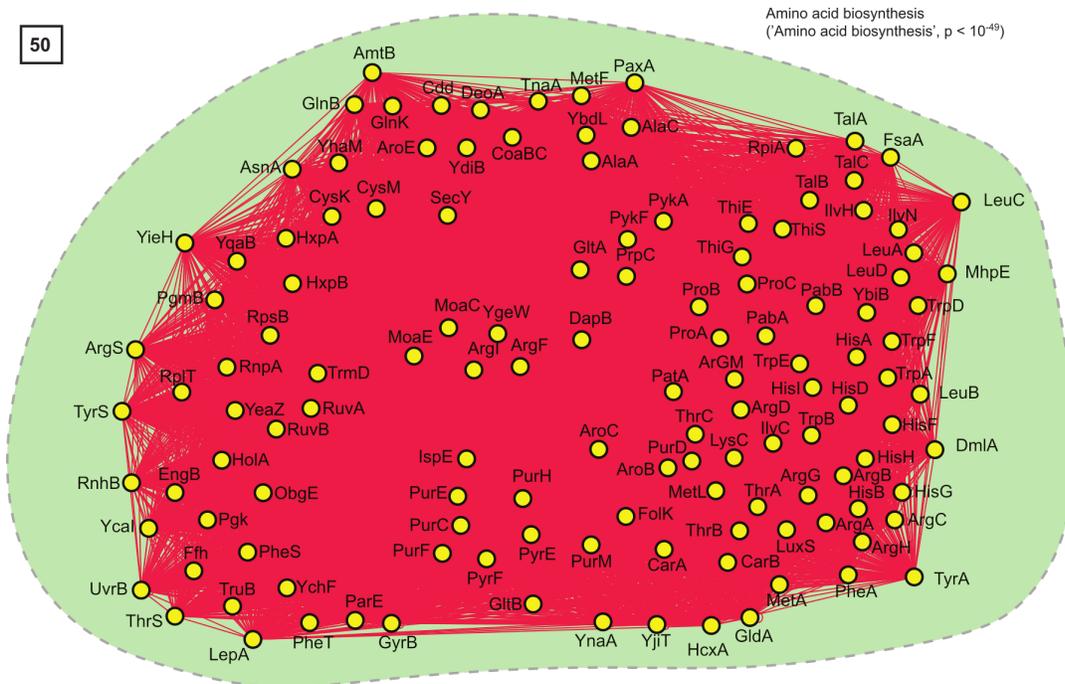
1612 identify discrete subnetworks and are labeled with their assigned function based on
1613 interpretation of gene-set enrichment analysis (GSEA) and literature review. The most
1614 significantly enriched ontological term produced by GSEA and the associated p-value is shown
1615 in parentheses for each subnetwork (**Table S9**). **B**) Comparison of the statistically-derived
1616 model using Hag (left) to the KEGG model (BR:bsu02035, right) of *B. subtilis* 168 motility. Venn
1617 diagrams represent the overlap between the sets of proteins in the indicated subnetwork of
1618 panel A (left) and the indicated KEGG category (right). **C**) Pie graph of the number of proteins
1619 in the statistical model that are represented in the KEGG hierarchy ('KEGG'), missing from the
1620 KEGG hierarchy but supported by experimental evidence in the literature ('Literature'), or absent
1621 from the KEGG hierarchy and the literature ('Novel').

1622

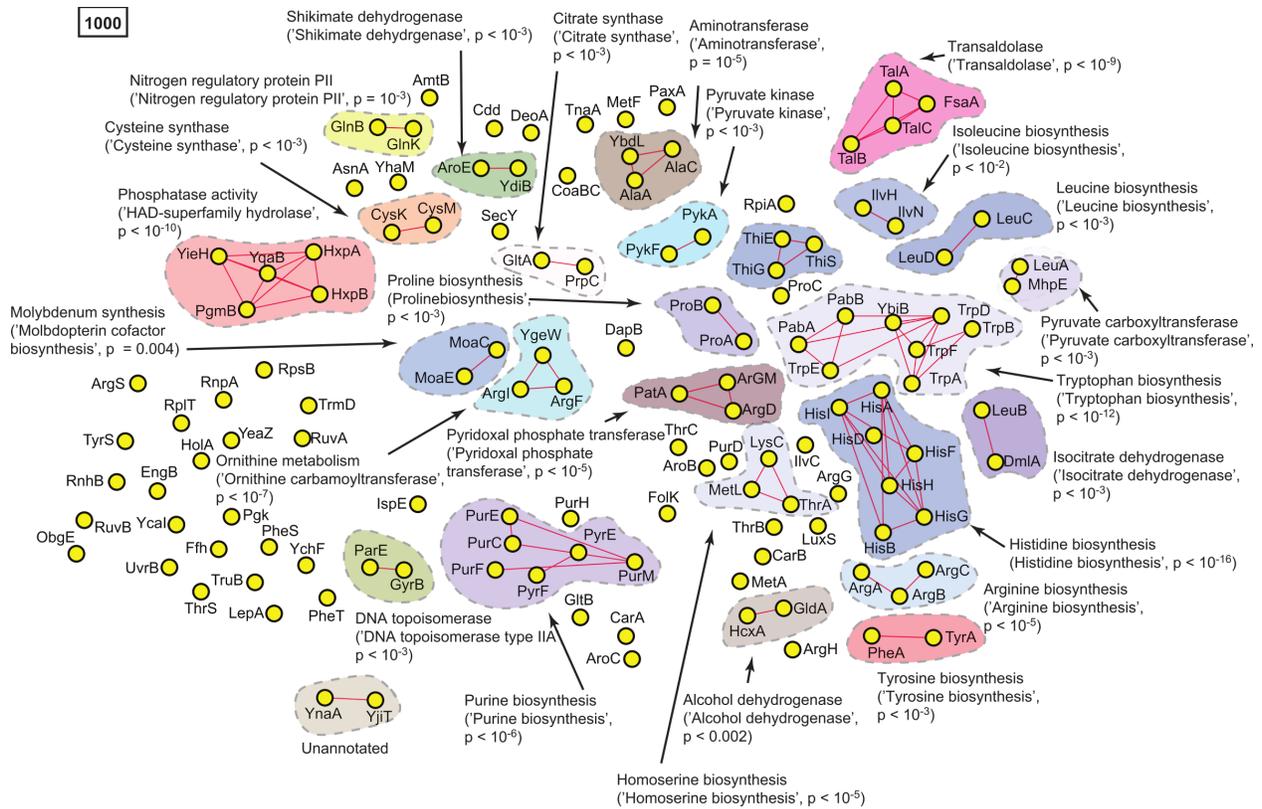
1623

1624

1625



1626



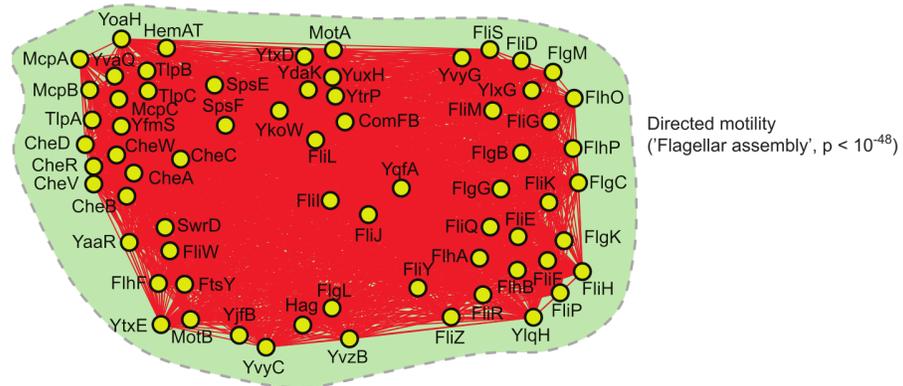
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640

Figure 4 – figure supplement 5

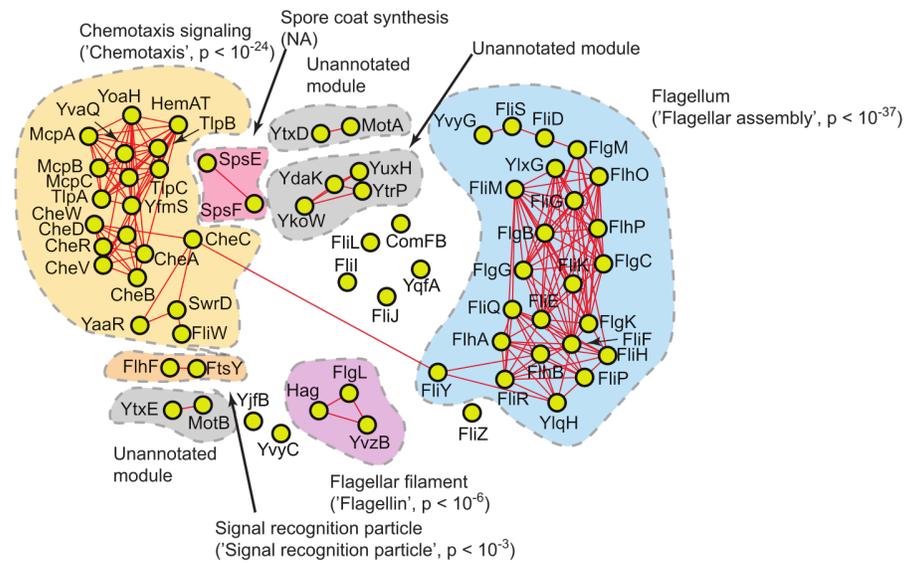
A statistically-derived hierarchical model of amino acid metabolism in *E. coli* K12 using HisG as a query protein. (A) 129 proteins in *E. coli* K12 shared significant spectral correlations with HisG across components SVD₃₄ to SVD₁₃₄. Statistical interaction networks were defined by thresholding spectral depth at 50, 150, 300, and 1000. Nodes (yellow circles) represent proteins and edges (red lines) represent statistical interactions. Shaded contours identify discrete subnetworks and are labeled with their assigned function based on interpretation of gene-set enrichment analysis (GSEA) and literature review. The most significantly enriched ontological term produced by GSEA and the associated p-value is shown in parentheses for each subnetwork (**Table S10**).

A.

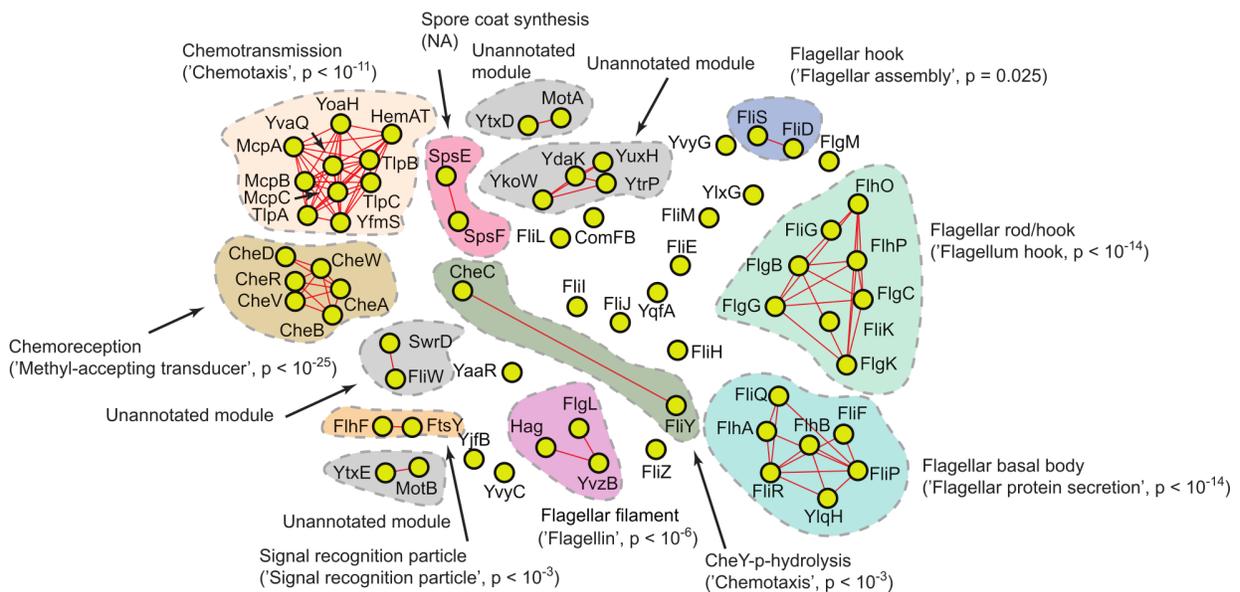
50



300



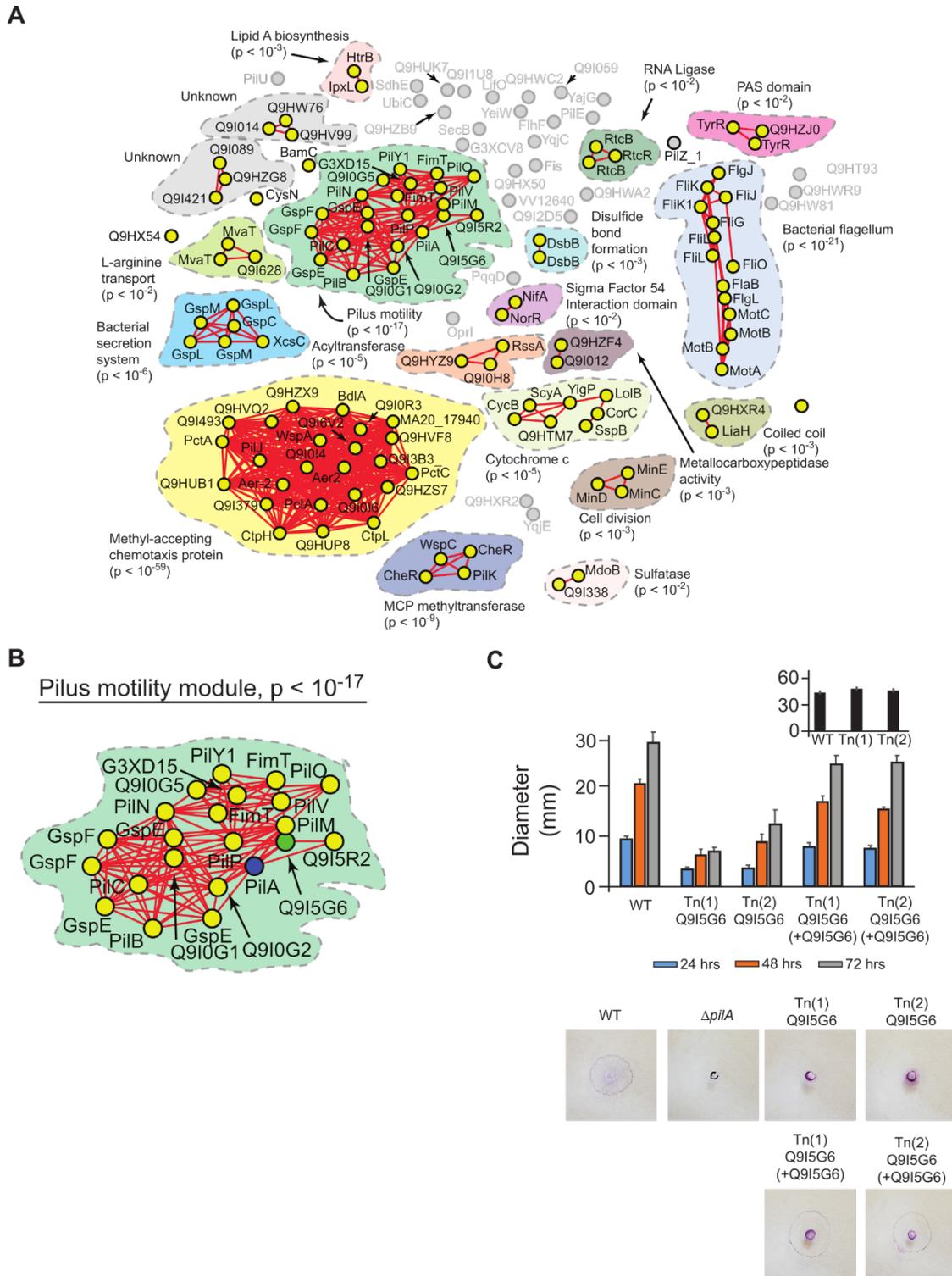
1000



1643 **Figure 4 – figure supplement 6**

1644

1645 **Statistically-derived hierarchical models of bacterial motility derived from serially**
1646 **thresholding spectral depth for correlations in *D^{domain}*. (A) A model of motility in *B.***
1647 ***subtilis* 168 using Hag as a query. (B) A model of motility in *E. coli* K12 using FliC as a**
1648 **query. Statistical interaction networks were defined by thresholding spectral depth at 50, 300,**
1649 **and 1000. Nodes (yellow circles) represent proteins and edges (red lines) represent statistical**
1650 **interactions. Shaded contours identify discrete subnetworks and are labeled with their assigned**
1651 **function based on interpretation of gene-set enrichment analysis (GSEA) and literature review.**
1652 **The most significantly enriched ontological term produced by GSEA and the associated p-value**
1653 **is shown in parentheses for each subnetwork (Table S11, S12).**



1654
1655

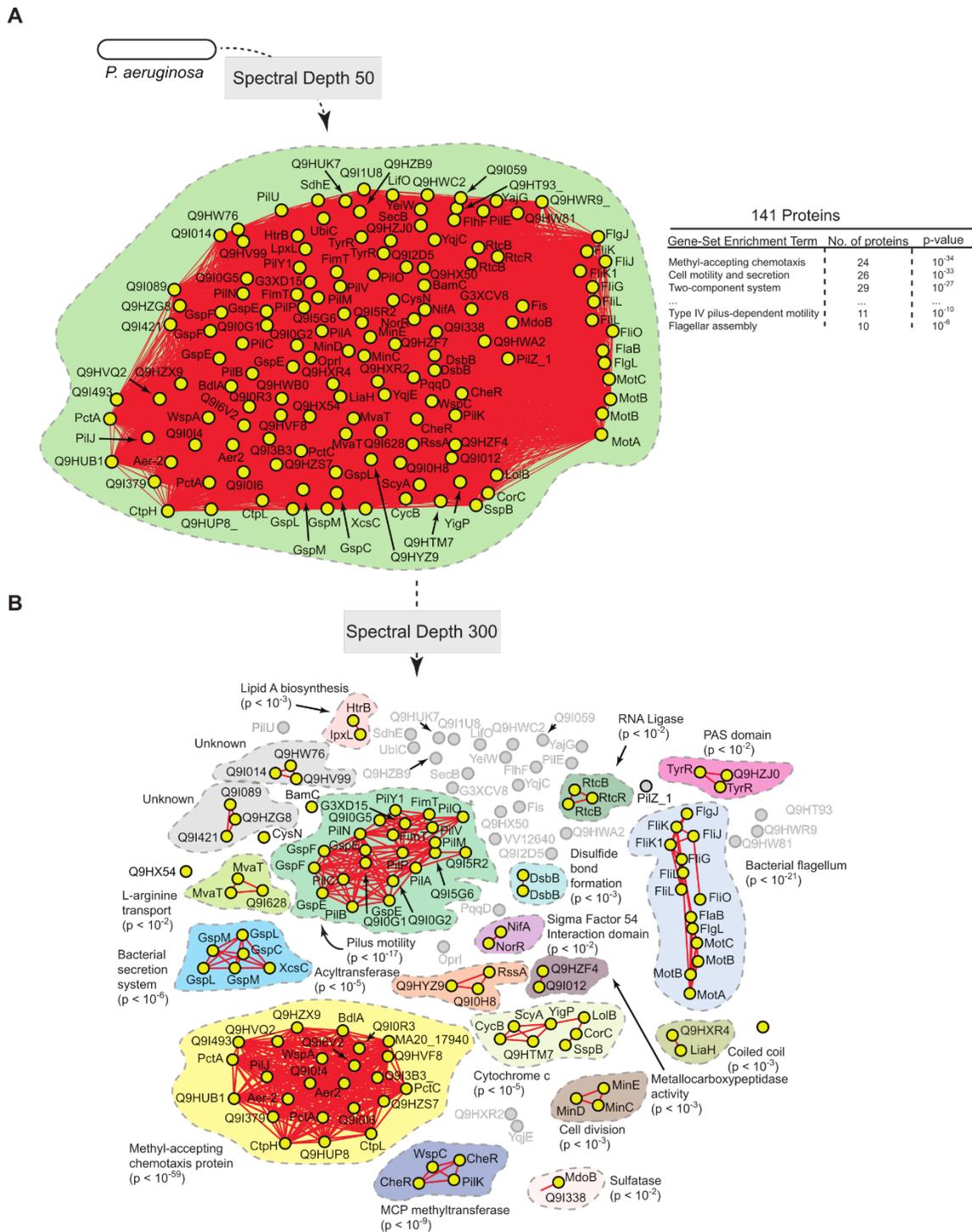
1656 **Figure 5. Prediction and experimental validation of a novel effector of twitch-based**
1657 **motility in *P. aeruginosa* (PAO1).** (A) Statistical network derived by applying a spectral depth
1658 threshold of 300 to the set of 141 protein in *P. aeruginosa* that were significantly correlated with
1659 PilA across SVD₃₄ to SVD₁₃₄. Nodes (yellow circles) represent proteins and edges (red lines)
1660 represent statistical interactions. Shaded contours identify statistical subnetworks and are

1661 labeled with their assigned function based on interpretation of gene-set enrichment analysis
1662 (GSEA) and literature review. The most significantly enriched ontological term produced by
1663 GSEA and the associated p-value is shown in parenthesis for each subnetwork (**Table S13**). (**B**)
1664 The pilus motility subnetwork from panel A. Nodes representing PilA and Q9I5G6 are colored
1665 blue and green respectively. (**C**) Time-course of pilus-based motility for parent (WT), two
1666 transposon mutants of Q9I5G6 (Tn(1) Q9I5G6, Tn(2) Q9I5G6), and transposon mutants
1667 complemented with Q9I5G6 (Tn(1) Q9I5G6 +Q9I5G6, Tn(2) Q9I5G6 +Q9I5G6). Inset shows
1668 results of flagellar motility for the parent strain (WT), and the two transposon mutants of Q9I5G6
1669 (Tn(1), Tn(2)) 24 hours post-inoculation. Representative images of the crystal-violet stained
1670 plates are shown.

1671

1672

1673



1674
1675
1676
1677
1678
1679
1680
1681
1682

Figure 5 – figure supplement 1

Statistically-derived hierarchical model of directed-motility in *P. aeruginosa* using PilA as a query. 140 proteins in *P. aeruginosa* shared significant spectral correlations with PilA across SVD₃₄ to SVD₁₃₄. **(A)** Statistical interaction network defined by thresholding spectral depth at 50. The inset illustrates significantly enriched terms resulting from gene-set enrichment analysis

1683 (GSEA) of the entire network. **(B)** Statistical interaction network defined by thresholding spectral
1684 depth at 300. Nodes (yellow circles) represent proteins and edges (red lines) represent
1685 statistical interactions. Shaded contours identify discrete subnetworks and are labeled in panel
1686 B with their assigned function based on interpretation of GSEA and literature review. The most
1687 significantly enriched ontological term produced by GSEA and the associated p-value is shown
1688 in parentheses for each subnetwork (**Table S13**).

1689

1690

1691

1692

1693

1694

1695