

## Decoding sequence-level information to predict membrane protein expression

Shyam M. Saladi<sup>1</sup>, Nauman Javed<sup>1</sup>, Axel Müller<sup>1</sup>, & William M. Clemons, Jr.<sup>1\*</sup>

<sup>1</sup>Division of Chemistry and Chemical Engineering, California Institute of Technology, Pasadena, CA

\*Corresponding author

Email: [clemons@caltech.edu](mailto:clemons@caltech.edu) (WMC)

## 1 Abstract

2 The expression of membrane proteins remains a major bottleneck in the characterization of these  
3 important proteins. Expression levels are currently unpredictable, which renders the pursuit of these  
4 targets challenging and inefficient. Evidence demonstrates that small changes in the nucleotide or  
5 amino-acid sequence can dramatically affect membrane protein biogenesis; yet these observations have  
6 not resulted in generalizable approaches to improve expression. Here, we develop a data-driven  
7 statistical model, named IMProve, that enriches for the likelihood of selecting membrane proteins that  
8 express in *E. coli* directly from sequence. The model, trained on experimental data, combines a set of  
9 sequence-derived variables resulting in a score that predicts the likelihood of expression. We test the  
10 model against various independent datasets that contain a variety of experimental outcomes  
11 demonstrating that the model significantly enriches for expressed proteins. Analysis of the underlying  
12 features reveals a significant role for nucleotide derived features in predicting expression. This  
13 computational model can immediately be used to identify favorable targets for characterization.

## 14 Author Summary

15 Membrane proteins play a pivotal role in biology, representing a quarter of all proteomes and a  
16 majority of drug targets. While considerable effort has been focused on improving our functional  
17 understanding of this class, much of the investment has been hampered by the inability to obtain  
18 sufficient amounts of sample. Until now, there have been no broadly successful strategies for predicting  
19 and improving expression which means that each target requires an *ad hoc* adventure. Complex  
20 biological processes govern membrane protein expression; therefore, sequence characteristics that  
21 influence protein biogenesis are not simply additive. Many properties must be considered  
22 simultaneously in predicting the expression level of a protein.

23 We provide a first solution to the membrane protein expression problem by learning from  
24 published data to develop a statistical model that predicts the outcomes of expression trials across  
25 families, scales, and laboratories (all independent of the model's training data). Given that the process of  
26 finding a target for large-scale expression is arduous, often requiring a long trial-and-error process that  
27 consumes significant financial and human resources, this work will have immediate applicability. The  
28 ability to study and engineer inaccessible membrane proteins becomes feasible with the use of our  
29 predictor. Furthermore, this work will enable others in developing new computational methods to assist  
30 in the experimental study of membrane proteins.

## 31 Introduction

32 The central role of integral membrane proteins motivates structural and biophysical studies that  
33 require large amounts of purified protein, often at considerable cost of both material and labor. Only a  
34 small percentage can be produced at high-levels resulting in membrane protein structural  
35 characterization lagging behind that of soluble proteins presently constituting just 1.7% of known  
36 atomic-level structures [1]. To increase the pace of structure determination, the scientific community  
37 created large government-funded structural genomics consortia facilities, like the NIH-funded New  
38 York Consortium on Membrane Protein Structure (NYCOMPS)[2]. For this representative example,  
39 more than 8000 genes, chosen based on characteristics hypothetically related to success, yielded only  
40 600 (7.1%) highly expressing proteins [3] resulting to date in 34 (5.6% of expressed proteins) unique

41 structures (based on annotation in the RCSB PDB [4]). This highlights the funnel problem of structural  
42 biology where each stage of the structure pipeline eliminates a large percentage of targets compounding  
43 into an overall low rate of success [5]. With new and rapidly advancing technologies like cryo-electron  
44 microscopy and micro-electron diffraction, we expect that the latter half of the funnel, structure  
45 determination, will increase in success rate [6,7]. In any case, membrane protein expression will  
46 continue to limit targets accessible for study [8].

47 Tools for improving the number of expressed membrane proteins are needed. While significant  
48 work has shown promise on a case-by-case basis, *e.g.* growth at lower temperatures, codon optimization  
49 [9], and regulating transcription [10], a generalizable solution remains elusive. Currently, each target  
50 must be addressed individually as the conditions that were successful for a previous target seldom carry  
51 over to other proteins, even amongst closely related homologs [5,11]. For individual cases, simple  
52 changes can have dramatic effects on the amount of expressed proteins [12,13]. Considering the  
53 scientific value of membrane protein studies, it is surprising that there are no methods that can provide  
54 solutions for improved expression outcomes with broad applicability across protein families and  
55 genomes.

56 Currently no approaches are available that decode sequence-level information for predicting  
57 membrane protein expression; yet the concept that sequence variation can measurably influence  
58 membrane protein biogenesis is commonplace. For example, positive-charges on cytoplasmic loops are  
59 important determinants of membrane protein topology [14,15]; yet introduction of mutations presumed  
60 to enhance certain properties, such as the positive inside rule, has not proven generalizable for  
61 improving expression [11]. The reasons for this likely lie in the complex underpinnings of membrane  
62 protein biogenesis, where the interplay between sequence features at the protein and nucleotide levels  
63 must be considered. Optimizing for a single sequence-level feature likely diminishes the beneficial  
64 effect of other features (*e.g.* increasing positive residues on internal loops might diminish favorable  
65 mRNA properties). Without accounting for the broad set of features related to membrane protein  
66 expression, it is impossible to predict differences in expression.

67 Attempts to develop algorithms that predict membrane protein expression have failed. Several  
68 examples, Daley, von Heijne, and coworkers [9,16,17] as well as NYCOMPS, were unable to use  
69 experimental expression data sets to train models that returned any predictive performance (personal  
70 communication). Statistical tools have been developed to predict expression and/or crystallization  
71 propensities from sequence information based on outcomes. These are primarily based on results from  
72 the Protein Structure Initiative where experimental outcomes are deposited in TargetTrack[18,19] and  
73 include well-known methods such as SPINE[20], Xtalpred[21–23], and PXS[24] as well as others[25–  
74 35]. While collectively these methods have supported significant advances in biochemistry, each suffers  
75 from similar issues when predicting membrane protein outcomes due to the criteria applied during the  
76 model development process. As membrane proteins have an extremely low success rate compared to  
77 soluble proteins, they are either explicitly excluded from the training process or are implicitly down-  
78 weighted by the statistical model. The result is that these methods do not predict membrane protein  
79 expression (representative methodology [21]).

80 In an ideal world, a perfect predictor would define the subset of protein sequences that can be  
81 expressed in a given host. As discussed elsewhere [9,16,17], none have successfully been able to map  
82 membrane protein expression to sequence. Given the scale of difficulty in expressing membrane  
83 proteins, we demonstrate here for the first time that it is possible to predict membrane protein expression  
84 purely based on sequence allowing one to enrich their expression trials for proteins with a higher  
85 probability of success.

86 To connect sequence to prediction, we develop a statistical model that maps a set of sequences to  
87 experimental expression levels via calculated features—thereby simultaneously accounting for the many  
88 potential determinants of expression. The resulting model allows ranking of any arbitrary set of  
89 membrane protein sequences in order of their relative likelihood of successful expression. In this first  
90 demonstration of prediction, we sought to select the simplest framework necessary to capture the  
91 problem. In particular, we train a linear equation that provides a score based on calculating the sum of  
92 weighted features where the weights are derived from fitting to experimental expression data, a “training  
93 set.” These features attempt to encapsulate the corpus of work that shows that sequence-level  
94 characteristics are important determinants of protein biogenesis, *e.g.* RNA secondary structure [36,37],  
95 transmembrane segment hydrophobicity [38–40], the positive inside rule [41], and loop disorder [42].

96 We extensively validate our model against a variety of independent datasets demonstrating its  
97 generalizability. This model can be used broadly to score any membrane protein based on its calculated  
98 features. In the process, we have built a method to enrich for positive expression outcomes with respect  
99 to the low positive rate attained from randomly selecting targets. To support further experimental efforts,  
100 we showcase the performance of the model across protein families and we broadly score the membrane  
101 proteome from a variety of important genomes. This approach and resulting model provides an exciting  
102 example for connecting sequence space to complex experimental outcomes.

## 103 Results

104 For this study, we focus on heterologous expression in *E. coli*, due to its ubiquitous use as a tool  
105 for membrane protein expression. While the benefits derived from low cost and low barriers for  
106 adoption are obvious, the applicability to the spectrum of the membrane proteome are becoming clearer.  
107 Of note, 43 of the 216 unique eukaryotic membrane protein structures were solved using protein  
108 expressed in *E. coli* (based on annotation in the RCSB PDB [4]). This demonstrates the utility of *E. coli*  
109 as a broad tool and its potential if the expression problem can be overcome.

### 110 Development of a computational model trained on *E. coli* expression data

111 A key component of any data-driven statistical model is the choice of dataset used for training.  
112 Having searched the literature, we identified two publications that contained quantitative datasets on the  
113 IPTG-induced overexpression of *E. coli* polytopic membrane proteins in *E. coli*. The first set, Daley,  
114 Rapp *et al.*, contained activity measures, proxies for expression level, from C-terminal tags of either  
115 GFP or PhoA (alkaline phosphatase)[16]. The second set, Fluman *et al.*, used a subset of constructs from  
116 the first and contained a more detailed analysis utilizing in-gel fluorescence to measure folded  
117 protein[43] (see Methods 4c). The expression results strongly correlated (Spearman’s  $\rho = 0.73$ ) between  
118 the two datasets demonstrating that normalized GFP activity was a good measure of the amount of  
119 folded membrane protein (Fig 1A and [43,44]). The experimental set-up employed multiple 96-well  
120 plates over multiple days resulting in pronounced variability in the absolute expression level of a given  
121 protein between trials. Daley, Rapp *et al.* calculated average expression levels by dividing the raw  
122 expression level of each protein by that of a control construct (Inverse LepB-GFP or LepB-PhoA) on the  
123 corresponding plate. While the resulting values were useful for the relevant question of identifying  
124 topology, we were unable to successfully fit a linear regression or a standard linear Support Vector  
125 Machine (SVM) to predict either the raw expression data compiled from all plates or averaged outcomes  
126 of each gene using numerical features calculated from nucleotide and protein sequences (see S1 Table;

127 Methods 2,3). This unexpected outcome suggested that the measurements required a more complex  
128 analysis.

129

130 **Fig 1. Training performance.** (A) A comparison of GFP activity [16] with measured folded protein  
131 [43] where each point represents the mean for a given gene tested in both works, and error bars plot the  
132 extrema. Spearman's rank correlation coefficient and 95% confidence interval (CI) [45] are shown. (B)  
133 Plates are the number of independent sets of measurements within which expression levels can be  
134 reliably compared. Genes are the number of proteins for which the C-terminus was reliably ascertained  
135 [16]. Observations are the total number of expression data points accessible. Total pairs are the number  
136 of comparable expression measurements (*i.e.* those within a single plate). Kendall's  $\tau$  is the metric  
137 maximized by the training process (See Methods 4b). The color of the column heading identifying each  
138 experimental set is retained throughout the figure. (C) Agreement against the normalized outcomes  
139 plotted as the mean activity (see Methods 5 for definition) versus the score with error bars providing the  
140 extent of observed activities (Spearman's  $\rho$  and 95% CI noted). (D) Illustrative Receiver Operating  
141 Characteristics (ROC) for thresholds at 25<sup>th</sup> and 75<sup>th</sup> percentile in activity with the number of positive  
142 outcomes at that threshold, the Area Under the Curve (AUC), and 95% CI indicated. (E) The AUC of  
143 the ROC at every possible activity threshold.

144

145 We hypothesized that measurements could be more accurately compared within an individual  
146 plate than across the entire dataset. To account for this, a preference-ranking linear SVM algorithm  
147 (SVM<sup>rank</sup> [46]) was chosen (see Methods 4b). Simply put, the SVM<sup>rank</sup> algorithm determines the optimal  
148 weight for each feature to best rank the order of expression outcomes within each plate over all plates,  
149 which results in a model where higher expressing proteins have higher scores. The outcome is identical  
150 in structure to a multiple linear regression, but instead of minimizing the sum of squared residuals, the  
151 SVM cost function is used accounting for the plate-wise constraint specified above. In practice, the  
152 process optimizes the correlation coefficient Kendall's  $\tau$ , as a training metric, to converge upon a set of  
153 weights. Kendall's  $\tau$  measures the agreement between ordinal quantities by calculating the number of  
154 correctly ordered and swapped pairs.

155 Various metrics related to the training data can be derived to assess the accuracy with which the  
156 model fits the input data (see Methods 4c). The SVM<sup>rank</sup> training metric shows varying agreement for all  
157 groups (*i.e.*,  $\tau_{\text{kendall}} > 0$ ) (Fig 1B). For individual genes, activity values normalized and averaged across  
158 trials were not directly used for the training procedure (see Methods 4a); yet one would anticipate that  
159 scores for each gene should broadly correlate with the expression average. Indeed, the observed  
160 normalized activities positively correlate with the score (dubbed IMProve score for Integral Membrane  
161 Protein expression improvement) output by the model (Fig 1C). Since SVM<sup>rank</sup> transforms raw  
162 expression levels within each plate to ranks before training, there is no expectation or guarantee that  
163 magnitude differences in expression level manifest in magnitude differences in score. As a result,  
164 Spearman's  $\rho$ , a rank correlation coefficient describing the agreement between two ranked quantities, is  
165 better suited for quantifying correlation over more common metrics like the  $R^2$  of a regression and  
166 Pearson's  $r$ .

167 For a more quantitative approach to assessing the model's success within the training data, we  
168 turn to the Receiver Operating Characteristic (ROC). ROC curves quantify the tradeoff between true  
169 positive and false positive predictions across the numerical scores output from a predictor. This is a  
170 more reliable assessment of prediction than simply calculating accuracy and precision from a single,



171 arbitrary score threshold [47]. The figure of merit that quantifies a ROC curve is the Area Under the  
172 Curve (AUC). Given that the AUC for a perfect predictor corresponds to 100% and that of a random  
173 predictor is 50% (Fig 1D, grey dashed line), an AUC greater than 50% indicates predictive performance  
174 of the model (percentage signs hereafter omitted) (see Methods 5 and [47]). Here, the ROC framework  
175 will be used to quantitatively assess the ability of our model to predict the outcomes within the various  
176 datasets.

177 The training datasets are quantitative measures of activity requiring that an activity threshold be  
178 chosen that defines positive or negative outcomes. For example, ROC curves using two distinct activity  
179 thresholds, at the 25<sup>th</sup> or 75<sup>th</sup> percentile of highest expression, are plotted with their calculated AUC  
180 values (Fig 1D). While both show that the model has predictive capacity, a more useful visualization  
181 would consider all possible activity thresholds. For this, the AUC value for every activity threshold is  
182 plotted showing that the model has predictive power regardless of an arbitrarily chosen expression  
183 threshold (Fig 1E). In total, the analysis demonstrates that the model can rank expression outcomes  
184 across all proteins in the training set. Interestingly, for PhoA-tagged proteins the model is progressively  
185 less successful with increasing activity. Since PhoA activity is an indirect measure of expression of  
186 proteins with their C-termini in the periplasm, this brings into question either the utility of this  
187 quantification method relative to GFP activity or perhaps that this class of proteins are special in the  
188 model. An argument for the former is presented later (Fig 2E).

189

190 **Fig 2. Success of the model against outcomes from NYCOMPS.** (A) An overview of the NYCOMPS  
191 outcomes and (B) a histogram of the number of conditions tested per gene colored based on outcome.  
192 (C) Receiver Operating Characteristics for positive groupings given by Only Positive outcomes genes  
193 (red) and genes with at least one positive outcome (pink). The percent positive for each group  
194 (corresponding color), total counts (black), and Area Under the Curve (AUC) values with 95%  
195 Confidence Interval (CI) are shown. The ROC considering genes with Mixed outcomes only as positive  
196 is shown as a blue dashed line with an AUC of 53.5 (51.8-55.2). The grey dashed line shows the  
197 performance of a completely random predictor (AUC = 50). (D) Histograms of genes with Only Positive  
198 (red) and Only Negative outcomes (grey) across IMProve scores (binned as described in Methods 5).  
199 The percentage of Only Positive outcomes in each bin is overlaid as a brown line (right axis). (E) The  
200 Positive Predictive Value (PPV) plotted for each percentile IMProve score, *e.g.* 75 on the x-axis  
201 indicates the PPV for the top 25% of genes based on score for genes, where positive indicates genes  
202 with Only Positive outcomes. The dashed line shows the overall success rate of the NYCOMPS  
203 experimental outcomes (~11% Only Positive). (F) The fold change in the PPV as a function of IMProve  
204 score relative to the success rate of NYCOMPS. (G) The AUCs for outcomes in each individual plasmid  
205 and solubilization condition (DDM except LDAO where noted) along with 95% CI (numerically in S2  
206 Table). Performances are also split by predicted C-terminal localization [48]. The numbers below  
207 indicate the total number of trials for each group and the percent within that group that were positive.

208

## 209 **Demonstration of prediction against an independent large expression dataset**

210 While the above analyses show that the model successfully fits the training data, we assess the  
211 broader applicability of the model outside the training set based on its success at predicting the outcomes  
212 of independent expression trials from distinct groups and across varying scales. The first test considers  
213 results from NYCOMPS, where 8444 membrane protein genes entered expression trials, in up to eight

214 conditions, resulting in 17114 expression outcomes (Fig 2A) [2]. The majority of genes were attempted  
215 in only one condition (Fig 2B), and, importantly, outcomes were non-quantitative (binary: expressed or  
216 not expressed) as indicated by the presence of a band by Coomassie staining of an SDS-PAGE gel after  
217 small-scale expression, solubilization, and nickel affinity purification [3]. For this analysis, the  
218 experimental results are either summarized as outcomes per gene or broken down as raw outcomes  
219 across defined expression conditions. For outcomes per gene, we can consider various thresholds for  
220 considering a gene as positive based on NYCOMPS expression success (Fig 2B). The most stringent  
221 threshold only regards a gene as positive if it has no negative outcomes (“Only Positive”, Fig 2B, red).  
222 Since a well expressing gene would generally advance in the NYCOMPS pipeline without further small-  
223 scale expression trials, this positive group likely contains the best expressing proteins. A second  
224 category comprises genes with at least one positive and at least one negative trial (“Mixed”, Fig 2B,  
225 blue). These genes likely include proteins that are more difficult to express.

226 ROCs assess predictive power across these groups (Fig 2C). IMProve scores markedly  
227 distinguish genes in the most stringent positive group (Only Positive) from all other genes (Fig 2C red).  
228 A permissive threshold considering genes as positive with at least one positive trial (Only Positive plus  
229 Mixed genes) shows more moderate predictive power (Fig 2C pink, AUC = 59.7 versus 67.1). If instead  
230 solely the Mixed genes are considered positive (excluding the Only Positive), the difference in the two  
231 positive groups is clear as the model very weakly distinguishes the mixed group from Only Negative  
232 genes (Fig 2C dashed blue, AUC = 53.5 (51.8-55.2)). This likely supports the notion that this pool  
233 largely consists of more difficult-to-express genes. For further analysis of NYCOMPS, we focus on the  
234 Only Positive pool as this likely represents the pool of best expressing proteins.

235 This predictive power can be qualitatively visualized as a histogram of the IMProve scores for  
236 genes separated by protein group (Only Positive, red; Only Negative, grey) (Fig 2D). Visually, the  
237 distribution of the scores for the Only Positive group is shifted to a higher score relative to the Only  
238 Negative group. This is emphasized considering the dramatic increase in the percentage of positive  
239 genes as a function of increasing IMProve score (overlaid as a brown line). A major aim of this work is  
240 to enrich the likelihood of choosing positively expressing proteins. The positive predictive value (PPV,  
241 true positives ÷ predicted positives) becomes a useful metric for positive enrichment as it conveys the  
242 degree of improved prediction over the experimental baseline of the dataset. The PPV of the model is  
243 plotted as a function of the percentile of the IMProve score for the Only Positive group (Fig 2E). In the  
244 figure, the experimental baseline is represented by a dashed line (11.1%); therefore, a relative increase  
245 reflects the predictive power of the algorithm. For example, considering the PPV of 20% for the top  
246 fourth of genes by IMProve score (75<sup>th</sup> percentile) shows that the algorithm increases the positive  
247 outcomes by 9% over baseline. For further illustration, we plot the fold-change in PPV across the  
248 various thresholds (Fig 2F). Here, if only genes with an IMProve score greater than -0.21 (75<sup>th</sup>  
249 percentile) were tested, the experiments would have returned nearly twice as many positives, a 1.82 fold  
250 change (Fig 2D). Higher score cut-offs would have even better returns.

251 Because there were eight different expression conditions, a final consideration looks at the  
252 NYCOMPS data based on the type of trial. Importantly, the model shows consistent performance  
253 throughout each of the eight conditions tested (Fig 2F, numerically in S2 Table). This highlights that the  
254 model is not sensitive to the experimental design of the training set and appears to predict broadly  
255 against different vector backbones. With this in mind, as an overall perspective, using a reasonable  
256 threshold for IMProve score (91<sup>st</sup> percentile or 0.5 (Fig 2E, yellow line)), had NYCOMPS tested the  
257 same number of genes an additional 1207 proteins would have been positive, representing a significant  
258 improvement in the return on investment.

259 The ability to predict the experimental data from NYCOMPS allows returning to the question of  
260 alkaline phosphatase as a metric for expression. For the training set, proteins with C-termini in the  
261 periplasm show less consistent fitting by the model (Fig 1, orange). To assess the generality of this  
262 result, the NYCOMPS outcomes are split into pools for either cytoplasmic or periplasmic C-terminal  
263 localization and AUCs are calculated for each. There are no significant differences in predictive capacity  
264 across all conditions (Fig 2G, green vs. orange) demonstrating that the model is applicable for all  
265 topologies.

### 266 **Further demonstration of prediction against small-scale independent datasets**

267 The NYCOMPS example demonstrates the predictive power of the model across the broad range  
268 of sequence space encompassed by that dataset. Next, the performance of the model is tested against  
269 relevant subsets of sequence space (*e.g.* a family of proteins or the proteome from a single organism),  
270 which are reminiscent of laboratory-scale experiments that precede structural or biochemical analyses.  
271 While a number of datasets exist [5,49–59], we identified six for which complete sequence information  
272 could be obtained to calculate all the necessary sequence features [49–54].

273 The first dataset is derived from the expression of 14 archaeal transporters in *E. coli* chosen  
274 based on their homology to human proteins [49]. For each putative transporter, expression was  
275 performed in three plasmids and two strains (six total conditions) with the membrane fraction quantified  
276 by both a Western blot against a histidine-affinity tag and Coomassie Blue staining of an SDS-PAGE  
277 gel. Here, the majority of the expressing proteins fall into the higher half of the IMProve scores, 7 out of  
278 9 of those with multiple positive outcomes (Fig 3A). Strikingly, quantification of the Coomassie Blue  
279 staining highlights a clear correlation with the IMProve score where the higher expressing proteins have  
280 the highest score (Fig 3B). ROC curves are plotted for the two thresholds: expression detected at least by  
281 Western blot or, for the smaller subset, by Coomassie Blue (Fig 3C). In both cases, the model shows  
282 predictive power. Consistent with what was seen for NYCOMPS, selecting only the top half of proteins  
283 by IMProve score would have captured the majority of the positive outcomes.

284

285 **Fig 3. Success of the model against a variety of small scale outcomes.** For each set, vertical lines  
286 indicate the median IMProve score. Receiver Operating Characteristics (ROC) along with Areas Under  
287 the Curves (AUC) and 95% confidence interval as well as the total number of positives for the given  
288 threshold (red hues) along with the total outcomes (black) are presented. In each curve, increasing  
289 expression thresholds as defined by the original publication are displayed as deeper red. **(A,B)** The  
290 expression of archaeal transporters in up to 6 trials. **(A)** Positive expression count is plotted above the  
291 dashed line and negative outcomes below the line. **(B)** From the same work, the expression of proteins  
292 detected by Coomassie Blue [49]. **(C)** ROC curves for each positive threshold (*i.e.* Coomassie Blue or  
293 Western Blot) from trials in **A,B**. **(D)** Experimental expression of *M. tuberculosis* membrane proteins  
294 plotted based on outcomes. **(E)** ROC curves for each possible threshold from trials in **D**. **(F)** Mammalian  
295 GPCR expression in either *E. coli* (top) or *P. pastoris* (bottom). **(G)** ROC curves for each possible  
296 threshold from trials in **F**.

297

298 The next test considers the expression of 105 *Mycobacterium tuberculosis* proteins in *E. coli*  
299 [50]. Protein expression was measured both by Coomassie Blue staining of an SDS-PAGE gel and  
300 Western blot with only outcomes from the membrane fraction considered for this analysis. The highest



301 expressing proteins (detected via Coomassie Blue) follow the trend given by the IMProve score with 7  
302 of the 9 falling within the higher half of scoring proteins (Fig 3D) and is reflected in the ROC (Figure  
303 3E). In contrast, using the positive Western blot outcomes as the minimum threshold (Fig 3D) shows an  
304 AUC no better than random (Fig 3E). Given that no internal standard was used and that each expression  
305 trial was performed only once, proteins that were positive by Western blot may represent a pool  
306 indistinguishable in expression from those not detected; alternatively, these results support that IMProve  
307 accurately captures the most highly expressing proteins. Again, selecting only the top half of the  
308 proteins based on their IMProve score would have captured nearly all of the high expressing proteins.

309 A broader test considers expression trials of 101 mammalian GPCRs in bacterial and eukaryotic  
310 systems [51]. Trials in *E. coli*, measured via Western blot of an insoluble fraction, again show highly  
311 expressing proteins at higher IMProve scores while the expression of the same proteins in *P. pastoris*,  
312 measured via dot blot, fail to show broad agreement (Fig 3F,G). The lack of predictive performance in  
313 *P. pastoris* suggests that the parameterization of the model, calibrated for broadly characterizing *E. coli*  
314 expression, requires retraining to generate a different model that captures the distinct interplay of  
315 sequence parameters in yeast. Still, the higher IMProve score clearly enriches for expressing proteins in  
316 *E. coli*.

317 Further expression trials of membrane proteins from *H. pylori*, *T. maritima* as well as microbial  
318 secondary transporters continues to show the same broad agreement [52–54] (S1 Fig). *H. pylori*  
319 membrane proteins showed that as the threshold for positive expressing proteins increases, the  
320 performance of the model improves (using the highest threshold  $n=46$  and  $AUC=67.7$ ) (S1 Fig. A,B).  
321 For *T. maritima* expression, the model weakly captures outcomes for two defined thresholds ( $n=5$  and  
322 19,  $AUC=61.7$  and 58.7), but due to the small number of successful outcomes, the confidence intervals  
323 are broad (S1 Fig. C,D). The expression of microbial secondary transporters shows varied agreement  
324 with the model. Taking proteins at the lower defined expression threshold shows predictive performance  
325 ( $n=59$ ,  $AUC=60.5$ ); however, considering the defined high-expressing proteins is less conclusive ( $n=26$ ,  
326  $AUC=52.0$ ) (S1 Fig. E,F). Broadly, independent of laboratory and experimental set-up, the IMProve  
327 score can enrich for the highest expressing proteins.

## 328 **Performance of the model across protein families**

329 To provide a clear path forward for experiment, we consider the performance of the model with  
330 regards to protein homology families, as defined by Pfam family classifications [60]. The 8444 genes in  
331 the NYCOMPS dataset fall into 555 families with ~15% not classified. To understand whether IMProve  
332 score is biased towards families present in the training set, we separate genes in the NYCOMPS dataset  
333 into three groups: part of the 153 families found in the training set, family not in the training set, and no  
334 defined Pfam family. There is no significant difference in AUC at 95% confidence between these groups  
335 (Fig 4A, bottom row). Therefore, the predictive power for a gene does not depend on the presence of its  
336 family within the training set.

337

338 **Fig 4. Model performance across protein families.** (A) The NYCOMPS dataset split by the presence  
339 or absence of a Pfam family in the training set with AUCs calculated by considering Only Positive genes  
340 as positive outcomes. (B) For each family within NYCOMPS with at least five outcomes (including one  
341 positive and one negative), the AUC across all outcomes is plotted with horizontal bars indicating the  
342 95% confidence interval. The color indicates the significance of the prediction within the family: purple,  
343 predictive at 95% confidence, blue, predictive but not at 95% confidence, green, not predictive. The size

344 of each significance group and total number of families (grey) are indicated on the plot. (C) Outcomes  
345 for specific protein families with an optimal IMProve score threshold indicated. Each was only tested in  
346 a single condition (N: His-FLAG-TEV-gene). CopD is classified as [TCDB 9.B.62](#) and AtoE as [TCDB](#)  
347 [2.A.73](#) [61]. (D) For the families in C, a ROC curve with the overall positive percentage within the  
348 group, total number of outcomes, and AUC with 95% CI is labelled.

349

350 The scale of NYCOMPS allows us to investigate whether there are protein families for which the  
351 model does better or worse than the aggregate. For this, an AUC is calculated for each protein family  
352 that has minimally five total outcomes (including at least one positive and one negative). Fig 4B plots  
353 the AUC for each protein family in increasing order as a cumulative distribution function. The breadth  
354 of the AUC values highlights the variability in predictive power across families. Most families can be  
355 predicted by the model (115 of 159 have an AUC > 0.5, visually blue and purple) though some not at  
356 95% confidence (57 of 115, blue), likely due to an insufficient number tested. Therefore, the  
357 NYCOMPS dataset provides some perspective on the protein families that IMProve best predicts.

358 For the protein families that are well-predicted within the NYCOMPS set, IMProve gives highly  
359 accurate insight into the likelihood of expression of a given protein. We demonstrate the utility of this  
360 prediction by looking at protein families that have yet to be characterized structurally. While there are a  
361 number of choices, one example is the protein family annotated as copper resistance proteins (CopD,  
362 PF05425), that typically contains eight transmembrane domains with an overall length of ~315 amino  
363 acids. A second example is the protein family annotated as short-chain fatty-acid transporters (AtoE,  
364 PF02667), that typically contains 10 transmembrane domains with an overall length of ~450 amino  
365 acids. In Fig 4C, genes from the two families are plotted by IMProve score and colored by outcome. In  
366 both cases, as indicated by the ROCs (Fig 4D), the model provides a clear score cut-off to guide target  
367 selection for future expression experiments. For example, considering CopD homologs, one would  
368 expect that those with IMProve scores above -1 will have a higher likelihood of expressing than on  
369 average across all homologs. This analysis can be broadly applied across the families that are predicted  
370 with high accuracy (S3 Table).

## 371 **Forward predictions on genomes of interest**

372 The model successfully enriches for heterologous expression of membrane proteins in *E. coli*  
373 strikingly across scales, laboratories, quantification methods, and protein families supporting its broad  
374 generalizability. While few genes express in every condition tested (Fig 2B and 3A), IMProve predicts  
375 the likelihood that a gene will express within a set of conditions and enriches for those that will work in  
376 any condition (Fig 2G, numerically in S2 Table).

377 To expand on the utility of this model, IMProve scores were calculated for membrane proteins  
378 from a variety of metazoan and microbial genomes (Fig 5A and S2 Fig. A). Many genomes have a  
379 significant proportion of proteins with high scores particularly evidenced by portions of the distributions  
380 ahead of the median in *E. coli* given by the vertical dashed line (Fig 5A). The likelihood for successful  
381 expression may be inferred by equating IMProve score with the PPV of Only Positive gene outcomes  
382 within the NYCOMPS dataset which rises significantly at scores above zero (Fig 5B). The range of  
383 scores spans those representative of high-expressing membrane proteins in both *E. coli* (Fig 1C) as well  
384 as in the NYCOMPS dataset (Fig 2C) and provides suggested targets for future biophysical studies (S4  
385 Table).

386

387 **Fig 5. Forward predictions of membrane protein expression for various genomes. (A)** Calculated  
388 scores for proteins from a variety of genomes (count in parentheses; complete set provided in S2 Fig. A)  
389 plotted as contours of kernel density estimates of the number of proteins at a given score. Amplitude is  
390 only relative within a genome. The dot indicates the median, and the lines depict quantities of an  
391 analogous Tukey boxplot[62,63]. The vertical line shows the median score in *E. coli* to provide context  
392 for other distributions. **(B)** PPV of Only Positive gene outcomes within the NYCOMPS dataset. **(C)**  
393 Distribution of overlap coefficients (see Methods 7) for each sequence parameter comparing the entire  
394 *E. coli* membrane proteome vs. the training set from *E. coli*. The dashed line provides a threshold  
395 separating the cluster of highly-related features from those with lower overlap. **(D-F)** A comparison of  
396 overlap coefficients with the training set between NYCOMPS and **(D)** all forward predictions (S2 Fig.  
397 A), **(E)** thermophilic genomes (orange), or **(F)** *P. falciparum*. Mean Absolute Deviation is indicated for  
398 each plot.

399

400 The predictions present several surprises at the biological level. One such is that the distribution  
401 of membrane proteins from representative thermophilic bacterial genomes have generally lower relative  
402 IMProve scores than other genomes, which implies that these proteins, on average, are harder to express  
403 in *E. coli*. This is in contrast to the many empirical examples of proteins from thermophiles which are  
404 often primary targets of biophysical characterization, although analysis of structural genomics data of  
405 soluble proteins suggests only a small crystallization advantage for this group [24]. In the case of the  
406 malarial parasite *P. falciparum*, the inverse trend is true with higher than anticipated relative IMProve  
407 scores despite the expectation that these proteins would be hard to express in *E. coli*. A possible cause  
408 for the distribution of scores may lie in the differences in the features that define the proteins in these  
409 particular groups. As the training set consists only of native *E. coli* sequences, the range of values for  
410 each feature in the training set may not represent the full range of possible values for the feature. For the  
411 special cases highlighted, perhaps the underlying sequence features fall into a poorly characterized  
412 subset of sequence space bringing into question the applicability of the model for these cases.

413 To address the utility of the model relative to differences in the sampling of sequence features,  
414 we measure the overlap of the distributions of sequence features used for prediction (S1 Table) for a  
415 given subset (see Methods 7) (S2 Fig B). Simply put, if two subsets contain the same distribution of  
416 sequence features the expectation is that a given feature should approach 100%. In the simplest case,  
417 comparing the distribution of sequences features in all *E. coli* membrane proteins against the subset used  
418 in the training set shows that the majority of features have overlap values over 75% (Fig 5C), which  
419 provides a lower threshold for similarity of sequence feature range. For NYCOMPS sequences, most of  
420 the overlap values relative to the training set are above the threshold. As this set shows predictive  
421 performance, comparison to the training set provides a baseline to assess the reliability of predictions  
422 within other subsets (Fig 5D-F, x-axis). In the first case (Fig 5D), there is a strong correlation between  
423 all the forward predictions and NYCOMPS, *i.e.* values are near the diagonal (quantified by a Mean  
424 Absolute Deviation (MAD) = 11.6), suggesting that differences in feature space do not significantly  
425 affect the predictive power of the model. For the thermophiles subset (Fig 5E), the values again are close  
426 to the diagonal (*i.e.* low MAD = 10.6) implying that the predictions are credible. *P. falciparum* (Fig 5F),  
427 on the other hand, clearly shows stark differences as most features fall below the 75% cut-off (MAD =  
428 29.0) bringing into question the reliability of these predictions. A training set with broader coverage of  
429 the feature space may generate a better predictor for all genomes.

## 430 **Biological importance of various sequence features**

431 Using a simple proof-of-concept linear model has allowed for a straightforward and useful  
432 predictor. Understanding if any single biological determinant is driving prediction may provide insight  
433 into membrane protein biogenesis and expression. With a linear model, as employed here, this task is  
434 ordinarily straightforward; assuming features are distributed identically and independently (“i.i.d.”), the  
435 weight assigned to each feature corresponds its relative importance. However, in our case, the input  
436 features do not satisfy these conditions, *i.e.* a lack of uniformity in feature distributions (S2 Fig B) and  
437 significant correlation between individual features (S3 Fig). As a result, during the training procedure,  
438 unequal weight is placed across correlating features that represent the same underlying biological  
439 phenomena, thereby, complicating the process of determining the biological underpinnings of the  
440 IMProve score. For example, the importance of transmembrane segment hydrophobicity is distributed  
441 between several features: among these the average  $\Delta G_{\text{insertion}}$  [40] of TM segments has a positive weight  
442 whereas average hydrophobicity, a correlating feature, has a negative weight (S1 Table, S3 Fig). As  
443 many features, such as those related to hydrophobicity, are correlated; conclusive information cannot be  
444 obtained simply using weights of individual features to interpret the relative importance of their  
445 underlying biological phenomena. We address this complication by coarsening our view of the features  
446 to two levels: First, we analyze features derived from protein versus those derived from nucleotide  
447 sequence, and then we look more closely at features groups after categorizing by biological phenomena.

448 The coarsest view of the features is a comparison of those derived from protein sequence versus  
449 those derived from nucleotide sequence. The summed weight for protein features is around zero,  
450 whereas for nucleotide features the summed weight is slightly positive suggesting that in comparison  
451 these features may be more important to the predictive performance of the model (Fig 6A). Within the  
452 training set, protein features more completely explain the score both via correlation coefficients (Fig 6B)  
453 as well as through ROC analysis (Fig 6C). However, comparison of the predictive performance of the  
454 two subsets of weights shows that the nucleotide features alone can give similar performance to the full  
455 model for the NYCOMPS dataset (Fig 6D). Within the small-scale datasets investigated, using only  
456 protein or nucleotide features shows no difference in predictive power at 95% confidence (Fig 6E). It is  
457 important to note that this does not suggest that protein features are not important for membrane protein  
458 expression. Instead, within the context of the trained model, nucleotide features are critical for predictive  
459 performance for a large and diverse dataset such as NYCOMPS. This finding corroborates growing  
460 literature that the nucleotide sequence holds significant determinants of biological processes [36,43,64–  
461 66].

462

463 **Fig 6. Feature contributions to the model.** (A) Classifying features by the type of sequence they are  
464 calculated from. (B) Considering the training set (as in Fig 1), Spearman correlation coefficients with  
465 95% confidence intervals using individual feature categories for each grouping of data within the  
466 training set of *E. coli* membrane proteins. Colors indicate the subset being assessed (green, whole cell  
467 GFP fluorescence; orange, alkaline phosphatase activity; purple, folded protein by in-gel fluorescence).  
468 (C) Protein/nucleotide feature dependence within the training set substantiated by the AUC of the ROC  
469 at every possible activity threshold for feature subsets independently (as in Fig 1E). (D) The AUC and  
470 95% confidence intervals using only protein or nucleotide features. (E) Protein/nucleotide feature  
471 dependence across small scale datasets shown as AUCs of the ROC along with 95% CI for the condition  
472 with the best overall predictive power (black).

473



474 To understand whether we may be able to provide more detailed evidence for feature  
475 importance, we collapse conceptually similar features into categories that allow for potential biological  
476 interpretation (S1 Table). As compared to the entire set of individual features, this process substantially  
477 reduces inter-feature correlation (S3 Fig, S4 Fig B). For example, the hydrophobicity group incorporates  
478 sequence features such as average hydrophobicity, maximum hydrophobicity,  $\Delta G_{\text{insertion}}$ , etc. The full list  
479 of groupings is provided in S1 Table and S3 Fig.

480 Analysis of categories suggests the phenomena that drive prediction. To visualize this, the  
481 collapsed weights are summarized in Fig 6B where each bar contains individual feature weights within a  
482 category. Features with a negative weight are stacked to the left of zero and those with a positive weight  
483 are stacked to the right. A red dot represents the sum of all weights, and the length of the bar gives the  
484 total absolute value of the combined weights within a category. Ranking the categories based on the sum  
485 of their weight suggests that some of categories play a more prominent role than others. These include  
486 properties related to transmembrane segments (hydrophobicity and TM size/count), codon pair score,  
487 loop length, and overall length/pI.

488 To explore the role of each category in prediction, the performance of the model is assessed  
489 using only features within a single category at a time. First understanding which categories perform well  
490 in the training set indicates which feature the model pulls information from and suggests hypotheses as  
491 to which categories ought to perform well across the validation datasets. Since the outcomes within the  
492 training set are real-valued, predictive power can be assessed via correlation coefficients with the  
493 predicted score yielding a single number (as in Fig 1C) or through AUCs across all possible expression  
494 thresholds (as in Fig 1D,E). Using the former metric, for simplicity, to assess the predictive capacity of  
495 feature subsets within the training set (Fig 6C) suggests several of interest with high correlation  
496 coefficients including 5' Codon Usage, Length/pI, Loop Length, and SD-like Sites. Only Length/pI  
497 shows some predictive across subsets of the NYCOMPS dataset (S4 Fig D).

498 Importantly, careful analysis of the training and large-scale testing dataset shows that no feature  
499 category independently drives the predictor. Excluding each individually does not significantly affect  
500 the overall predictive performance, except for Length/pI (isoelectric point) (S4 Fig D). Sequence length  
501 composes the majority of the weight within this category and is one of the highest weighted features in  
502 the model. This is consistent with the anecdotal observation that larger membrane proteins are typically  
503 harder to express. However, this parameter alone would not be useful for predicting within a smaller  
504 subset, like a single protein family, where there is little variance in length (*e.g.* Fig 3,4). One might  
505 develop a predictor that was better for a given protein family under certain conditions with a subset of  
506 the entire features considered here; yet this would require *a priori* knowledge of the system, *i.e.* which  
507 sequence features were truly most important, and would preclude broad generalizability as shown for the  
508 predictor presented here.

## 509 **Sequence optimization for expression**

510 The predictive performance of the model implies that the features defined here provide a coarse  
511 approximation of the fitness landscape for membrane protein expression. Attempting to optimize a  
512 single feature by modifying the sequence will likely affect the resulting score and expression due to  
513 changes in other features. Fluman, *et al.* provides an illustrative experiment [43]. They hypothesized that  
514 altering the number of Shine-Dalgarno (SD)-like sites in the coding sequence of a membrane protein  
515 would affect expression. To test this, silent mutations were engineered within the first 200 bases of three  
516 proteins (genes *ygdD*, *brnQ*, and *ybjJ* from *E. coli*) to increase the number of SD-like sites with the goal  
517 of improving expression. Expression trials demonstrated that only one of the proteins (BrnQ) had  
518 improved expression of folded protein (Fig 7). However, the resulting changes in the IMProve score



519 correspond with the changes in measured expression as the model considers changes to other nucleotide  
520 features. Capture of the outcomes in this small test case by the model illustrates the utility of integrating  
521 the contribution of the numerous parameters involved in membrane protein biogenesis.

522

523 **Fig 7. Synonymous mutations affect expression.** Relative difference in SD-like sites (green),  
524 expression (purple), and IMProve score (yellow) between wild-type and mutants with silent mutations  
525 engineered to increase anti-SD sequence binding propensity [43]. See Methods 7 for further detail.

526

## 527 Discussion

528 Here, we have demonstrated the ability to predict membrane protein expression using  
529 computational methods, a feat some have considered impossible. Our success is built on encompassing a  
530 multitude of experimental results into a single computational model. The predictive power of IMProve  
531 provides a low barrier-to-entry method to enrich for positive expression outcomes.

532 The current best practice for characterization of a membrane protein target begins with the  
533 identification and testing of many homologs or variants for expression. IMProve will allow for  
534 prioritization of targets to test for expression thereby making more optimal use of limited human and  
535 material resources. In addition, due to the scale of NYCOMPS, protein families that were extensively  
536 tested provide ranges of scores (*e.g.* Fig 5C) where the score of an individual target directly indicates its  
537 likelihood of expression relative to known experimental results. We provide the current predictor as web  
538 service where scores can be calculated, and the method, associated data, and suggested analyses are  
539 publically available to catalyze progress across the community ([clemonslab.caltech.edu](http://clemonslab.caltech.edu)).

540 Having shown that membrane protein expression can be predicted, the generalizability of the  
541 model is remarkable despite several known limitations. Using data from a single study for training  
542 precludes including certain variables that empirically influence expression such as the features  
543 corresponding to fusion tags and the context of the protein in an expression plasmid, *e.g.* the 5'  
544 untranslated region, for which there was no variation in the Daley, Rapp, *et al.* dataset. Moreover, using  
545 a simple proof-of-concept linear model allowed for a straightforward and robust predictor; however,  
546 intrinsically it cannot be directly related to the biological underpinnings. While we can extract some  
547 biological inference, a linear combination of sequence features does not explicitly reflect the reality of  
548 physical limits for host cells. To some extent, constraint information is likely encoded in the complex  
549 architecture of the underlying sequence space (*e.g.* through the genetic code, TM prediction, RNA  
550 secondary structure analyses). Future statistical models that improve on these limitations will likely hone  
551 predictive power and more intricately characterize the interplay of variables that underlie membrane  
552 protein expression in *E. coli* and other systems.

553 A perhaps surprising outcome of our results is the demonstration of the quantitatively important  
554 contribution of the nucleotide sequence as a component of the IMProve score. This echoes the growing  
555 literature that aspects of the nucleotide sequence are important determinants of protein biogenesis in  
556 general [36,43,64–66]. While one expects that there may be different weights for various nucleotide  
557 derived features between soluble and membrane proteins, it is likely that these features are important for  
558 soluble proteins as well. An example of this is the importance of codon optimization for soluble protein  
559 expression, which has failed to show any general benefit for membrane proteins [9]. Current expression  
560 predictors that have predictive power for soluble proteins have only used protein sequence for deriving

561 the underlying feature set [22,35]. Future prediction methods will likely benefit from including  
562 nucleotide sequence features as done here.

563 The ability to predict phenotypic results using sequence based statistical models opens a variety  
564 of opportunities. As done here, this requires a careful understanding of the system and its underlying  
565 biological processes enumerated in a multitude of individual variables that impact the stated goal of the  
566 predictor, in this case enriching protein expression. As new features related to expression are discovered,  
567 future work will incorporate these leading to improved models. Based on these results, expanding to  
568 new expression hosts such as eukaryotes seems entirely feasible, although a number of new features may  
569 need to be considered, *e.g.* glycosylation sites and trafficking signals. Moreover, the ability to score  
570 proteins for expressibility creates new avenues to computationally engineer membrane proteins for  
571 expression. The proof-of-concept described here required significant work to compile data from  
572 genomics consortia and the literature in a readily useable form. As data becomes more easily accessible,  
573 broadly leveraging diverse experimental outcomes to decode sequence-level information, an extension  
574 of this work, is anticipated.

## 575 **Methods**

576 Sequence mapping & retrieval and feature calculation was performed in Python 2.7 [67] using  
577 BioPython [68] and NumPy [69]; executed and consolidated using Bash (shell) scripts; and parallelized  
578 where possible using GNU Parallel [70]. Data analysis and presentation was done in R [71] within  
579 RStudio [72] using magrittr [73], plyr [74], dplyr [75], asbio [76], and datamart [77] for data handling;  
580 ggplot2 [78], ggbeeswarm [79], GGally [80], gridExtra [81], cowplot [82], scales [83], viridis [84], and  
581 RColorBrewer [85,86] for plotting; multidplyr [87] with parallel [71] and foreach [88] with iterators  
582 [89] and doMC [90]/doParallel [91] for parallel processing; and roxygen2 [92] for code organization and  
583 documentation as well as other packages as referenced.

### 584 585 **1. Collection of data necessary for learning and evaluation**

586 ***E. coli* Sequence Data** – The nucleotide sequences from [16] were deduced by reconstructing forward  
587 and reverse primers (*i.e.* ~20 nucleotide stretches) from each gene in Colibri (based on EcoGene 11), the  
588 original source cited and later verified these primers against an archival spreadsheet provided directly by  
589 Daniel Daley (personal communication). To account for sequence and annotation corrections made to  
590 the genome after Daley, Rapp, *et al.*'s work, these primers were directly used to reconstruct the  
591 amplified product from the most recent release of the *E. coli* K-12 substr. MG1655 genome [93]  
592 (EcoGene 3.0; U00096.3). Although Daniel Daley mentioned that raw reads from the Sanger sequencing  
593 runs may be available within his own archives, it was decided that the additional labor to retrieve this  
594 data and parse these reads would not significantly impact the model. The deduced nucleotide sequences  
595 were verified against the protein lengths given in S1 Table from [16]. The plasmid library tested in [43]  
596 was provided by Daniel Daley, and those sequences are taken to be the same.

597  
598 ***E. coli* Training Data** – The preliminary results using the mean-normalized activities echoed the  
599 findings of [16] that these do not correlate with sequence features either in the univariate sense (many  
600 simple linear regressions, S1 Table [16]) or a multivariate sense (multiple linear regression, data not  
601 shown). This is presumably due to the loss of information regarding variability in expression level for  
602 given genes or due to the increase in variance of the normalized quantity (See Methods 4a) due to the  
603 normalization and averaging procedure. Daniel Daley and Mikaela Rapp provided spreadsheets of the  
604 outcomes from the 96-well plates used for their expression trials and sent scanned copies of the readouts  
605 from archival laboratory notebooks where the digital data was no longer accessible (personal  
606 communication). Those proteins without a reliable C-terminal localization (as given in the original  
607 work) or without raw expression outcomes were not included in further analyses.

608 Similarly, Nir Fluman also provided spreadsheets of the raw data from the set of three expression  
609 trials performed in [43].

610  
611 **New York Consortium on Membrane Protein Structure (NYCOMPS) Data** – Brian Kloss, Marco  
612 Punta, and Edda Kloppman provided a dataset of actions performed by the NYCOMPS center including  
613 expression outcomes in various conditions [2,3]. The protein sequences were mapped to NCBI GenInfo  
614 Identifier (GI) numbers either via the Entrez system [94] or the Uniprot mapping service[95]. Each GI  
615 number was mapped to its nucleotide sequence via a combination of the NCBI Elink mapping service  
616 and the “coded\_by” or “locus” tags of Coding Sequence (CDS) features within GenBank entries.  
617 Though a custom script was created, a script from Peter Cock on the BioPython listserv to do the same  
618 task via a similar mapping mechanism was found [96]. To confirm all the sequences, the TargetTrack  
619 [18] XML file was parsed for the internal NYCOMPS identifiers and compared for sequence identity to

620 those that had been mapped using the custom script; 20 (less than 1%) of the sequences had minor  
621 inconsistencies and were manually replaced.

622

623 **Archaeal transporters Data** – The locus tags (“Gene Name” in Table 1) were mapped directly to the  
624 sequences and retrieved from NCBI [49]. Pikyee Ma and Margarida Archer clarified questions regarding  
625 their work to inform the analysis.

626

627 **GPCR Expression Data** – Nucleotide sequences were collected by mapping the protein identifiers  
628 given in Table 1 from [51] to protein GIs via the Uniprot mapping service [95] and subsequently to their  
629 nucleotide sequences via the custom mapping script described above (see NYCOMPS). The sequence  
630 length and pI were validated against those provided. Renaud Wagner assisted in providing the  
631 nucleotide sequences for genes whose listed identifiers were unable to be mapped and/or did not pass the  
632 validation criteria as the MeProtDB (the sponsor of the GPCR project) does not provide a public  
633 archive.

634

635 ***Helicobacter pylori* Data** – Nucleotide sequences were retrieved by mapping the locus tags given in  
636 Supplemental Table 1 from [52] to locus tags in the Jan 31, 2014 release of the *H. pylori* 26695 genome  
637 (AE000511.1). To verify sequence accuracy, sequences whose molecular weight matched that given by  
638 the authors were accepted. Those that did not match, in addition to the one locus tag that could not be  
639 mapped to the Jan 31, 2014 genome version, were retrieved from the Apr 9, 2015 release of the genome  
640 (NC\_000915.1). Both releases are derived from the original sequencing project [97]. After this curation,  
641 all mapped sequences matched the reported molecular weight.

642

In this data set, expression tests were performed in three expression vectors and scored as 1, 2, or  
643 3. Two vectors were scored via two methods. For these two vectors, the two scores were averaged to  
644 give a single number for the condition making them comparable to the third vector while yielding 2  
645 additional thresholds (1.5 and 2.5) result in the 5 total curves shown (S1 Fig. B).

646

647 ***Mycobacterium tuberculosis* Data** – The authors note using TubercuList through GenoList [98],  
648 therefore, nucleotide sequences were retrieved from the archival website based on the original  
649 sequencing project [99]. The sequences corresponding to the identifiers and outcomes in Table 1 from  
650 [50] were validated against the provided molecular weight .

651

652 **Secondary Transporter Data** – GI Numbers given in Table 1 from [54] were matched to their CDS  
653 entries using the custom mapping script described above (see NYCOMPS). Only expression in *E. coli*  
654 with IPTG-inducible vectors was considered.

655

656 ***Thermotoga maritima* Data** – Gene names given in Table 1 [100] were matched to CDS entries in the  
657 Jan 31, 2014 release of the *Thermotoga maritima* MSB8 genome (AE000512.1), a revised annotation of  
658 the original release[101]. The sequence length and molecular weight were validated against those  
659 provided.

660

## 661 **2. Calculation of sequence features**

662

Based on experimental analyses and anecdotal evidence, approximately 105 different protein and  
663 nucleotide sequence features thought to be relevant to expression were identified and calculated for each  
664 protein using custom code together with published software (codonW [102], tAI [103], NUPACK [104],  
665 Vienna RNA [105], Codon Pair Bias [106], Disembl [42], and RONN [107]). Relative metrics (*e.g.*

666 codon adaptation index) are calculated with respect to the *E. coli* K-12 substr. MG1655 [93] quantity.  
667 The octanol-water partitioning [39], GES hydrophobicity [38],  $\Delta G$  of insertion [40] scales were  
668 employed as well. Transmembrane segment topology was predicted using Phobius Constrained for the  
669 training data and Phobius for all other datasets [48]. We were able to obtain the Phobius code and  
670 integrate it directly into our feature calculation pipeline resulting in significantly faster speeds than any  
671 other option. Two RNA secondary structure metrics were prompted in part by Goodman, et al. [36].  
672 Several features were obtained by averaging per-site metrics (*e.g.* per-residue RONN3.2 disorder  
673 predictions) in windows of a specified length. Windowed tAI metrics are calculated over *all* 30 base  
674 windows (not solely over 10 codon windows). S1 Table lists a description of each feature. Features are  
675 calculated solely from a gene of interest excluding portions of the ORFs such as linkers and tags derived  
676 from the plasmid backbone employed (future work will explore contributions of these elements).

### 677 678 **3. Preparation for model learning**

679 Calculated sequence features for the membrane proteins in the *E. coli* dataset as well as raw  
680 activity measurements, *i.e.* each 96-well plate, were loaded into R. As is best practice in using Support  
681 Vector Machines, each feature was “centered” and “scaled” where the mean value of a given feature was  
682 subtracted from each data point and then divided by the standard deviation of that feature using  
683 `preprocess` [108]. As is standard practice, the resulting set was then culled for those features of near  
684 zero-variance, over 95% correlation (Pearson’s  $r$ ), and linear dependence (`nearZeroVar`,  
685 `findCorrelation`, `findLinearCombos`)[108]. In particular this procedure removed extraneous  
686 degrees of freedom during the training process which carry little to no additional information with  
687 respect to the feature space and which may over represent certain redundant features. Features and  
688 outcomes for each list (“query”) were written into the SVM<sup>light</sup> format using a modified  
689 `svmlight.write` [109].

690 The final features were calculated for each sequence in the test datasets, prepared for scoring by  
691 “centering” and “scaling” by the training set parameters via `preprocess` [108], and then written into  
692 SVM<sup>light</sup> format again using a modified `svmlight.write`.

### 693 694 **4. Model selection, training, and evaluation using SVM<sup>rank</sup>**

695 **a.** At the most basic level, our predictive model is a learned function that maps the parameter space  
696 (consisting of nucleotide and protein sequence features) to a response variable (expression level)  
697 through a set of governing weights ( $w_1, w_2, \dots, w_N$ ). Depending on how the response variable is defined,  
698 these weights can be approximated using several different methods. As such, defining a response  
699 variable that is reflective of the available training data is key to selecting an appropriate learning  
700 algorithm.

701 The quantitative 96-well plate results [16] that comprise our training data do not offer an  
702 absolute expression metric valid over all plates—the top expressing proteins in one plate would not  
703 necessarily be the best expressing within another. As such, this problem is suited for preference-ranking  
704 methods. As a ranking problem, the response variable is the ordinal rank for each protein derived from  
705 its overexpression relative to the other members of the same plate of expression trials. In other words,  
706 the aim is to rank highly expressed proteins (based on numerous trials) at higher scores than lower  
707 expressed proteins by fitting against the order of expression outcomes from each constituent 96-well  
708 plate.

709 **b.** As the first work of this kind, the aim was to employ the simplest framework necessary taking in  
710 account the considerations above. The method chosen computes all valid pairwise classifications (*i.e.*  
711 within a single plate) transforming the original ranking problem into a binary classification problem.



712 The algorithm outputs a score for each input by minimizing the number of swapped pairs thereby  
713 maximizing Kendall's  $\tau$  [110]. For example, consider the following data generated via context A  
714  $(X_{A,1}, Y_{A,1}), (X_{A,2}, Y_{A,2})$  and B  $(X_{B,1}, Y_{B,1}), (X_{B,2}, Y_{B,2})$  where observed response follows as index  $i$ , *i.e.*  
715  $Y_n < Y_{n+1}$ . Binary classifier  $f(X_i, X_j)$  gives a score of 1 if an input pair matches its ordering criteria and  
716  $-1$  if not, *i.e.*  $Y_i < Y_j$ :

$$\begin{aligned} 717 \quad & f(X_{A,1}, X_{A,2}) = 1; f(X_{A,2}, X_{A,1}) = -1 \\ 718 \quad & f(X_{B,1}, X_{B,2}) = 1; f(X_{B,2}, X_{B,1}) = -1 \\ 719 \quad & f(X_{A,1}, X_{B,2}), f(X_{A,2}, X_{B,1}) \text{ are invalid} \end{aligned}$$

720 Free parameters describing  $f$  are calculated such that those calculated orderings  
721  $f(X_{A,1}), f(X_{A,2}) \dots; f(X_{B,1}), f(X_{B,2}) \dots$  most closely agree (overall Kendall's  $\tau$ ) with the observed  
722 ordering  $Y_n, Y_{n+1}, \dots$ . In this sense,  $f$  is a pairwise Learning to Rank method.

723 Within this class of models, a linear preference-ranking Support Vector Machine was employed  
724 [111]. To be clear, as an algorithm a preference-ranking SVM operates similarly to the canonical SVM  
725 binary classifier. In the traditional binary classification problem, a linear SVM seeks the maximally  
726 separating hyper-plane in the feature space between two classes, where class membership is determined  
727 by which side of the hyper-plane points reside. For some  $n$  linear separable training examples  $D =$   
728  $\{(x_i) | x_i \in \mathbb{R}^d\}^n$  and two classes  $y_i \in \{-1, 1\}$ , a linear SVM seeks a mapping from the  $d$ -dimensional  
729 feature space  $\mathbb{R}^d \rightarrow \{-1, 1\}$  by finding two maximally separated hyperplanes  $w \cdot x - b = 1$  and  $w \cdot$   
730  $x - b = -1$  with constraints that  $w \cdot x_i - b \geq 1$  for all  $x_i$  with  $y_i \in \{1\}$  and  $w \cdot x_i - b \leq -1$  for all  
731  $x_i$  with  $y_i \in \{-1\}$ . The feature weights correspond to the vector  $w$ , which is the vector perpendicular to  
732 the separating hyperplanes, and are computable in  $O(n \log n)$  implemented as part of the SVM<sup>rank</sup>  
733 software package, though in  $O(n^2)$  [46]. See [111] for an in-depth, technical discussion.

734 **c.** In a soft-margin SVM where training data is not linearly separable, a tradeoff between misclassified  
735 inputs and separation from the hyperplane must be specified. This parameter  $C$  was found by training  
736 models against raw data from Daley, Rapp, *et al.* with a grid of candidate  $C$  values ( $2^n \forall n \in [-5, 5]$ )  
737 and then evaluated against the raw "folded protein" measurements from Fluman, *et al.* The final model  
738 was chosen by selecting that with the lowest error from the process above ( $C = 2^5$ ). To be clear, the final  
739 model is composed solely of a single weight for each feature; the tradeoff parameter  $C$  is only part of the  
740 training process.

741 Qualitatively, such a preference-ranking method constructs a model that ranks groups of proteins  
742 with higher expression level higher than other groups with lower expression value. In comparison to  
743 methods such as linear regression and binary classification, this approach is more robust and less  
744 affected by the inherent stochasticity of the training data.

## 745 5. Quantitative Assessment of Predictive Performance

746 In generating a predictive model, one aims to enrich for positive outcomes while ensuring they  
747 do not come at the cost of increased false positive diagnoses. This is formalized in Receiver Operating  
748 Characteristic (ROC) theory (for a primer see [47]), where the true positive rate is plotted against the  
749 false positive rate for all classification thresholds (score cutoffs in the ranked list). In this framework, the  
750 overall ability of the model to resolve positive from negative outcomes is evaluated by analyzing the  
751 Area Under a ROC curve (AUC) where  $AUC_{\text{perfect}} = 100\%$  and  $AUC_{\text{random}} = 50\%$  (percentage signs are  
752 omitted throughout the text and figures). All ROCs are calculated through pROC [112] using the  
753 analytic Delong method for AUC confidence intervals [113]. Bootstrapped AUC CIs ( $N = 10^6$ ) were  
754 precise to 4 decimal places suggesting that analytic CIs are valid for the NYCOMPS dataset.

756 With several of our datasets, no definitive standard or clear-cut classification for positive  
757 expression exists. However, the aim is to show and test all reasonable classification thresholds of  
758 positive expression for each dataset in order to evaluate predictive performance as follows:

759 **Training data** – The outcomes are quantitative (activity level), so each ROC is calculated by  
760 normalizing within each dataset to the standard well subject to the discussion in 4a above (LepB for  
761 PhoA, and InvLepB for GFP) (examples in Fig 1D) for each possible threshold, *i.e.* each normalized  
762 expression value with each AUC plotted in Fig 1E. 95% confidence intervals of Spearman's  $\rho$  are given  
763 by  $10^6$  iterations of a bias-corrected and accelerated (BCa) bootstrap of the data (Fig 1A,C) [45].

764 **Large-scale** – ROCs were calculated for each of the expression classes (Fig 2E). Regardless of the split,  
765 predictive performance is noted. The binwidth for the histogram was determined using the Freedman-  
766 Diaconis rule[114], and scores outside the plotted range comprising  $<0.6\%$  of the density were implicitly  
767 hidden.

768 **Small-scale** – Classes can be defined in many different ways. To be principled about the matter, ROCs  
769 for each possible cutoff are presented based on definitions from each publication (Fig 3C,E,G, S1 Fig.  
770 B,D,F). See Methods 1 for any necessary details about outcome classifications for each dataset.

771

## 772 **6. Feature Weights**

773 Weights for the learned SVM are pulled directly from the model file produced by SVM<sup>light</sup> and are given  
774 in S1 Table.

775

## 776 **7. Forward Predictions**

777 **Data collection** – We selected several genomes for comparison as shown in Fig 5, S2 Fig. A, and S4  
778 Table. Coding sequences of membrane proteins from human and mouse genomes were gathered by  
779 mapping Uniprot identifiers of proteins noted to have at least one transmembrane segment by Uniprot  
780 [95] to Ensembl (release 82) coding sequences [115] via Biomart [116]. *C. elegans* coding sequences  
781 were similarly mapped via Uniprot but to WormBase coding sequences [117] also via Biomart. *S.*  
782 *cerevisiae* strain S288C coding sequences [118] were retrieved from the Saccharomyces Genome  
783 Database. *P. pastoris* strain GS115 coding sequences [119] were retrieved from the DOE Joint Genome  
784 Institute (JGI) Genome Portal [120]. Those sequences without predicted [48] TMs were excluded from  
785 subsequent analyses. Microbial sequences were gathered via a custom, in-house database populated with  
786 data compiled primarily from Pfam [60], DOE JGI Integrated Microbial Genomes [121], and the  
787 Microbial Genome Database [122].

788 **Feature calculation** – Because of the incredible number of sequences, we did not calculate the features  
789 derived from the most computationally expensive calculation (whole sequence mRNA pairing  
790 probability). Since predictive performance on the NYCOMPS dataset is slightly smaller, but not  
791 significantly different at 95% confidence, in the absence of these features (S2 Table), the forward  
792 predictions are still valid. For future experiments, these features can be calculated for the subset of  
793 targets of interest.

794 **Parameter space similarity** – As a first approximation of the similarity of the  $\sim 90$  dimensional  
795 sequence parameter space between two groupings, features were compared pairwise via the following  
796 metric. Let  $f_i$  and  $g_i$  represent the true distributions for a given feature  $i$  between two groups of interest.  
797 The distribution overlap, *i.e.* shared area,  $\Delta_i$  is formalized as

798 
$$\Delta_i(f_i, g_i) = \int \min\{f_i(x), g_i(x)\} dx$$

799 ranging from 0, for entirely distinct distributions, to 1 for entirely identical distributions.

800 As written  $f_i$  and  $g_i$  are probability densities, they need to be approximated before calculating  $\Delta_i$   
801 and are done so via kernel density estimates (KDE) of the observed samples  $[x_1^f, \dots, x_n^f]$  and  $[x_1^g, \dots, x_n^g]$   
802 using a nonparametric, locally adaptive method allowing for variable bandwidth smoothing  
803 implemented in LocFit[123] ( $\text{adpen}=2\sigma^2$ ) providing  $\hat{f}_i$  and  $\hat{g}_i$ . The distribution overlap  $\Delta_i$  is evaluated  
804 over a grid of  $2^{13}$  equally spaced points over the range of  $f_i$  and  $g_i$ .

805 *Shine-Dalgarno-like mutagenesis* – Folded protein is quantified by densitometry measurement [124,125]  
806 of the relevant band in Figure 6 of [43]. Relative difference is calculated as is standard:

$$\frac{\text{metric}_{\text{mutant}} - \text{metric}_{\text{wildtype}}}{\frac{1}{2} |\text{metric}_{\text{mutant}} - \text{metric}_{\text{wildtype}}|}$$

809

## 8. Availability

810 All analysis is documented in a series of R notebooks[126] available openly at  
811 [github.com/lemlab/IMProve](https://github.com/lemlab/IMProve). These notebooks provide fully executable instructions for the  
812 reproduction of the analyses and the generation of figures and statistics in this study. The ranking engine  
813 is available as a web service at [clemonslab.caltech.edu](https://clemonslab.caltech.edu). Additional code is available upon request.

814

815

## Acknowledgements

816 We thank Daniel Daley and Thomas Miller's group for discussion, Yaser Abu-Mostafa and  
817 Yisong Yue for guidance regarding machine learning, Niles Pierce for providing NUPACK source code  
818 [104], Welison Floriano and Naveed Near-Ansari for maintaining local computing resources, and  
819 Samuel Schulte for suggesting the model's name. We thank James Bowie, Michiel Niesen, Stephen  
820 Marshall, Thomas Miller, Reid van Lehn, and Tom Rapoport for critical reading of the manuscript.  
821 Models and analyses are possible thanks to raw experimental data provided by Daniel Daley and  
822 Mikaela Rapp [16]; Nir Fluman [43]; Edda Kloppmann, Brian Kloss, and Marco Punta from  
823 NYCOMPS [2,3]; Pikyee Ma [49]; Renaud Wagner [51]; and Florent Bernaudat [55].

824 Computational time was provided by Stephen Mayo and Douglas Rees. This work used the  
825 Extreme Science and Engineering Discovery Environment (XSEDE), which is supported by National  
826 Science Foundation grant number ACI-1053575 [127].  
827

## 828 **References**

- 829 1. Hendrickson WA. Atomic-level analysis of membrane-protein structure. *Nat Struct Mol Biol.* 2016;23: 464–467.  
830 doi:10.1038/nsmb.3215
- 831 2. Punta M, Love J, Handelmann S, Hunt JF, Shapiro L, Hendrickson WA, et al. Structural genomics target selection for the  
832 New York consortium on membrane protein structure. *J Struct Funct Genomics.* 2009;10: 255–268.  
833 doi:10.1007/s10969-009-9071-1
- 834 3. Love J, Mancina F, Shapiro L, Punta M, Rost B, Girvin M, et al. The New York Consortium on Membrane Protein  
835 Structure (NYCOMPS): a high-throughput platform for structural genomics of integral membrane proteins. *J Struct*  
836 *Funct Genomics.* 2010;11: 191–199. doi:10.1007/s10969-010-9094-7
- 837 4. Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, et al. The Protein Data Bank. *Nucleic Acids Res.*  
838 2000;28: 235–242.
- 839 5. Lewinson O, Lee AT, Rees DC. The funnel approach to the precrystallization production of membrane proteins. *J Mol*  
840 *Biol.* 2008;377: 62–73. doi:10.1016/j.jmb.2007.12.059
- 841 6. Merk A, Bartesaghi A, Banerjee S, Falconieri V, Rao P, Davis MI, et al. Breaking Cryo-EM Resolution Barriers to  
842 Facilitate Drug Discovery. *Cell.* 2016;165: 1698–1707. doi:10.1016/j.cell.2016.05.040
- 843 7. Nannenga BL, Gonen T. MicroED opens a new era for biological structure determination. *Curr Opin Struct Biol.*  
844 2016;40: 128–135. doi:10.1016/j.sbi.2016.09.007
- 845 8. Bill RM, Henderson PJF, Iwata S, Kunji ERS, Michel H, Neutze R, et al. Overcoming barriers to membrane protein  
846 structure determination. *Nat Biotechnol.* 2011;29: 335–340. doi:10.1038/nbt.1833
- 847 9. Nørholm MHH, Light S, Virkki MTI, Elofsson A, von Heijne G, Daley DO. Manipulating the genetic code for  
848 membrane protein production: what have we learnt so far? *Biochim Biophys Acta.* 2012;1818: 1091–1096.  
849 doi:10.1016/j.bbamem.2011.08.018
- 850 10. Wagner S, Klepsch MM, Schlegel S, Appel A, Draheim R, Tarry M, et al. Tuning *Escherichia coli* for membrane protein  
851 overexpression. *Proc Natl Acad Sci U S A.* 2008;105: 14371–14376. doi:10.1073/pnas.0804090105
- 852 11. Marshall SS, Niesen MJM, Müller A, Tiemann K, Saladi SM, Galimidi RP, et al. A Link between Integral Membrane  
853 Protein Expression and Simulated Integration Efficiency. *Cell Rep.* 2016;16: 2169–2177.  
854 doi:10.1016/j.celrep.2016.07.042
- 855 12. Sarkar CA, Dodevski I, Kenig M, Dudli S, Mohr A, Hermans E, et al. Directed evolution of a G protein-coupled  
856 receptor for expression, stability, and binding selectivity. *Proc Natl Acad Sci U S A.* 2008;105: 14808–14813.  
857 doi:10.1073/pnas.0803103105
- 858 13. Schlinkmann KM, Honegger A, Türeci E, Robison KE, Lipovšek D, Plückthun A. Critical features for biosynthesis,  
859 stability, and functionality of a G protein-coupled receptor uncovered by all-versus-all mutations. *Proc Natl Acad Sci*  
860 *U S A.* 2012;109: 9810–9815. doi:10.1073/pnas.1202107109
- 861 14. Seppälä S, Slusky JS, Lloris-Garcerá P, Rapp M, von Heijne G. Control of membrane protein topology by a single C-  
862 terminal residue. *Science.* 2010;328: 1698–1700. doi:10.1126/science.1188950
- 863 15. Van Lehn RC, Zhang B, Miller TF. Regulation of multispinning membrane protein topology via post-translational  
864 annealing. *eLife.* 2015;4. doi:10.7554/eLife.08697

- 865 16. Daley DO, Rapp M, Granseth E, Melén K, Drew D, von Heijne G. Global topology analysis of the Escherichia coli inner  
866 membrane proteome. *Science*. 2005;308: 1321–1323. doi:10.1126/science.1109730
- 867 17. Nørholm MHH, Toddo S, Virkki MTI, Light S, von Heijne G, Daley DO. Improved production of membrane proteins in  
868 Escherichia coli by selective codon substitutions. *FEBS Lett*. 2013;587: 2352–2358.  
869 doi:10.1016/j.febslet.2013.05.063
- 870 18. Chen L, Oughtred R, Berman HM, Westbrook J. TargetDB: a target registration database for structural genomics  
871 projects. *Bioinformatics*. 2004;20: 2860–2862. doi:10.1093/bioinformatics/bth300
- 872 19. Gabanyi MJ, Adams PD, Arnold K, Bordoli L, Carter LG, Flippen-Andersen J, et al. The Structural Biology  
873 Knowledgebase: a portal to protein structures, sequences, functions, and methods. *J Struct Funct Genomics*. 2011;12:  
874 45–54. doi:10.1007/s10969-011-9106-2
- 875 20. Bertone P, Kluger Y, Lan N, Zheng D, Christendat D, Yee A, et al. SPINE: an integrated tracking database and data  
876 mining approach for identifying feasible targets in high-throughput structural proteomics. *Nucleic Acids Res*.  
877 2001;29: 2884–2898.
- 878 21. Slabinski L, Jaroszewski L, Rodrigues APC, Rychlewski L, Wilson IA, Lesley SA, et al. The challenge of protein  
879 structure determination--lessons from structural genomics. *Protein Sci Publ Protein Soc*. 2007;16: 2472–2482.  
880 doi:10.1110/ps.073037907
- 881 22. Slabinski L, Jaroszewski L, Rychlewski L, Wilson IA, Lesley SA, Godzik A. XtalPred: a web server for prediction of  
882 protein crystallizability. *Bioinforma Oxf Engl*. 2007;23: 3403–3405. doi:10.1093/bioinformatics/btm477
- 883 23. Jahandideh S, Jaroszewski L, Godzik A. Improving the chances of successful protein structure determination with a  
884 random forest classifier. *Acta Crystallogr D Biol Crystallogr*. 2014;70: 627–635. doi:10.1107/S1399004713032070
- 885 24. Price WN, Chen Y, Handelman SK, Neely H, Manor P, Karlin R, et al. Understanding the physical properties that  
886 control protein crystallization by analysis of large-scale experimental data. *Nat Biotechnol*. 2009;27: 51–57.  
887 doi:10.1038/nbt.1514
- 888 25. Overton IM, Padovani G, Girolami MA, Barton GJ. ParCrys: a Parzen window density estimation approach to protein  
889 crystallization propensity prediction. *Bioinforma Oxf Engl*. 2008;24: 901–907. doi:10.1093/bioinformatics/btn055
- 890 26. Mizianty MJ, Kurgan L. Meta prediction of protein crystallization propensity. *Biochem Biophys Res Commun*.  
891 2009;390: 10–15. doi:10.1016/j.bbrc.2009.09.036
- 892 27. Babnigg G, Joachimiak A. Predicting protein crystallization propensity from protein sequence. *J Struct Funct Genomics*.  
893 2010;11: 71–80. doi:10.1007/s10969-010-9080-0
- 894 28. Kandaswamy KK, Pugalenti G, Suganthan PN, Gangal R. SVMCRYST: an SVM approach for the prediction of protein  
895 crystallization propensity from protein sequence. *Protein Pept Lett*. 2010;17: 423–430.
- 896 29. Mizianty MJ, Kurgan L. Sequence-based prediction of protein crystallization, purification and production propensity.  
897 *Bioinforma Oxf Engl*. 2011;27: i24–33. doi:10.1093/bioinformatics/btr229
- 898 30. Jahandideh S, Mahdavi A. RFCRYST: sequence-based protein crystallization propensity prediction by means of random  
899 forest. *J Theor Biol*. 2012;306: 115–119. doi:10.1016/j.jtbi.2012.04.028
- 900 31. Mizianty MJ, Kurgan LA. CRYSpred: accurate sequence-based protein crystallization propensity prediction using  
901 sequence-derived structural characteristics. *Protein Pept Lett*. 2012;19: 40–49.



- 902 32. Fusco D, Barnum TJ, Bruno AE, Luft JR, Snell EH, Mukherjee S, et al. Statistical analysis of crystallization database  
903 links protein physico-chemical features with crystallization mechanisms. *PLoS One*. 2014;9: e101123.  
904 doi:10.1371/journal.pone.0101123
- 905 33. Gao J, Hu G, Wu Z, Ruan J, Shen S, Hanlon M, et al. Improved Prediction of Protein Crystallization, Purification and  
906 Production Propensity Using Hybrid Sequence Representation. *Curr Bioinforma*. 2014;9: 57–64.  
907 doi:10.2174/15748936113080990006
- 908 34. Hu J, Han K, Li Y, Yang J-Y, Shen H-B, Yu D-J. TargetCrys: protein crystallization prediction by fusing multi-view  
909 features with two-layered SVM. *Amino Acids*. 2016;48: 2533–2547. doi:10.1007/s00726-016-2274-4
- 910 35. Wang H, Feng L, Zhang Z, Webb GI, Lin D, Song J. CrysAlis: an integrated server for computational analysis and design  
911 of protein crystallization. *Sci Rep*. 2016;6: 21383. doi:10.1038/srep21383
- 912 36. Goodman DB, Church GM, Kosuri S. Causes and effects of N-terminal codon bias in bacterial genes. *Science*.  
913 2013;342: 475–479. doi:10.1126/science.1241934
- 914 37. Mirzadeh K, Martínez V, Toddo S, Guntur S, Herrgård MJ, Elofsson A, et al. Enhanced Protein Production in  
915 *Escherichia coli* by Optimization of Cloning Scars at the Vector-Coding Sequence Junction. *ACS Synth Biol*. 2015;  
916 doi:10.1021/acssynbio.5b00033
- 917 38. Engelman DM, Steitz TA, Goldman A. Identifying nonpolar transbilayer helices in amino acid sequences of membrane  
918 proteins. *Annu Rev Biophys Biophys Chem*. 1986;15: 321–353. doi:10.1146/annurev.bb.15.060186.001541
- 919 39. Wimley WC, Creamer TP, White SH. Solvation energies of amino acid side chains and backbone in a family of host-  
920 guest pentapeptides. *Biochemistry (Mosc)*. 1996;35: 5109–5124. doi:10.1021/bi9600153
- 921 40. Hessa T, Meindl-Beinker NM, Bernsel A, Kim H, Sato Y, Lerch-Bader M, et al. Molecular code for transmembrane-  
922 helix recognition by the Sec61 translocon. *Nature*. 2007;450: 1026–1030. doi:10.1038/nature06387
- 923 41. Heijne G. The distribution of positively charged residues in bacterial inner membrane proteins correlates with the trans-  
924 membrane topology. *EMBO J*. 1986;5: 3021–3027.
- 925 42. Linding R, Jensen LJ, Diella F, Bork P, Gibson TJ, Russell RB. Protein disorder prediction: implications for structural  
926 proteomics. *Structure*. 2003;11: 1453–1459.
- 927 43. Fluman N, Navon S, Bibi E, Pilpel Y. mRNA-programmed translation pauses in the targeting of *E. coli* membrane  
928 proteins. *eLife*. 2014;3. doi:10.7554/eLife.03440
- 929 44. Geertsma ER, Groeneveld M, Slotboom D-J, Poolman B. Quality control of overexpressed membrane proteins. *Proc*  
930 *Natl Acad Sci U S A*. 2008;105: 5722–5727. doi:10.1073/pnas.0802190105
- 931 45. Cauty A, Ripley BD. *boot: Bootstrap R (S-Plus) Functions*. 2015.
- 932 46. Tsochantaridis I, Joachims T, Hofmann T, Altun Y. Large Margin Methods for Structured and Interdependent Output  
933 Variables. *J Mach Learn Res*. 2005;6: 1453–1484.
- 934 47. Swets JA, Dawes RM, Monahan J. Better decisions through science. *Sci Am*. 2000;283: 82–87.
- 935 48. Käll L, Krogh A, Sonnhammer ELL. A combined transmembrane topology and signal peptide prediction method. *J Mol*  
936 *Biol*. 2004;338: 1027–1036. doi:10.1016/j.jmb.2004.03.016
- 937 49. Ma P, Varela F, Magoch M, Silva AR, Rosário AL, Brito J, et al. An efficient strategy for small-scale screening and  
938 production of archaeal membrane transport proteins in *Escherichia coli*. *PLoS One*. 2013;8: e76913.  
939 doi:10.1371/journal.pone.0076913

- 940 50. Korepanova A, Gao FP, Hua Y, Qin H, Nakamoto RK, Cross TA. Cloning and expression of multiple integral  
941 membrane proteins from *Mycobacterium tuberculosis* in *Escherichia coli*. *Protein Sci.* 2005;14: 148–158.  
942 doi:10.1110/ps.041022305
- 943 51. Lundstrom K, Wagner R, Reinhart C, Desmyter A, Cherouati N, Magnin T, et al. Structural genomics on membrane  
944 proteins: comparison of more than 100 GPCRs in 3 expression systems. *J Struct Funct Genomics.* 2006;7: 77–91.  
945 doi:10.1007/s10969-006-9011-2
- 946 52. Psakis G, Nitschkowski S, Holz C, Kress D, Maestre-Reyna M, Polaczek J, et al. Expression screening of integral  
947 membrane proteins from *Helicobacter pylori* 26695. *Protein Sci.* 2007;16: 2667–2676. doi:10.1110/ps.073104707
- 948 53. Dobrovetsky E, Lu ML, Andorn-Broza R, Khutoreskaya G, Bray JE, Savchenko A, et al. High-throughput production of  
949 prokaryotic membrane proteins. *J Struct Funct Genomics.* 2005;6: 33–50. doi:10.1007/s10969-005-1363-5
- 950 54. Surade S, Klein M, Stolt-Bergner PC, Muenke C, Roy A, Michel H. Comparative analysis and “expression space”  
951 coverage of the production of prokaryotic membrane proteins for structural genomics. *Protein Sci.* 2006;15: 2178–  
952 2189. doi:10.1110/ps.062312706
- 953 55. Bernaudat F, Frelet-Barrand A, Pochon N, Dementin S, Hivin P, Boutigny S, et al. Heterologous expression of  
954 membrane proteins: choosing the appropriate host. *PloS One.* 2011;6: e29191. doi:10.1371/journal.pone.0029191
- 955 56. Eshaghi S, Hedrén M, Nasser MIA, Hammarberg T, Thornell A, Nordlund P. An efficient strategy for high-throughput  
956 expression screening of recombinant integral membrane proteins. *Protein Sci.* 2005;14: 676–683.  
957 doi:10.1110/ps.041127005
- 958 57. Gordon E, Horsefield R, Swarts HGP, de Pont JJHM, Neutze R, Snijder A. Effective high-throughput overproduction  
959 of membrane proteins in *Escherichia coli*. *Protein Expr Purif.* 2008;62: 1–8. doi:10.1016/j.pep.2008.07.005
- 960 58. Petrovskaya LE, Shulga AA, Bocharova OV, Ermolyuk YS, Kryukova EA, Chupin VV, et al. Expression of G-protein  
961 coupled receptors in *Escherichia coli* for structural studies. *Biochem Mosc.* 2010;75: 881–891.
- 962 59. Szakonyi G, Leng D, Ma P, Bettaney KE, Saidijam M, Ward A, et al. A genomic strategy for cloning, expressing and  
963 purifying efflux proteins of the major facilitator superfamily. *J Antimicrob Chemother.* 2007;59: 1265–1270.  
964 doi:10.1093/jac/dkm036
- 965 60. Finn RD, Bateman A, Clements J, Coghill P, Eberhardt RY, Eddy SR, et al. Pfam: the protein families database. *Nucleic  
966 Acids Res.* 2014;42: D222–230. doi:10.1093/nar/gkt1223
- 967 61. Saier MH, Reddy VS, Tsu BV, Ahmed MS, Li C, Moreno-Hagelsieb G. The Transporter Classification Database  
968 (TCDB): recent advances. *Nucleic Acids Res.* 2016;44: D372–379. doi:10.1093/nar/gkv1103
- 969 62. Tukey JW. *Exploratory data analysis.* Reading, Mass: Addison-Wesley Pub. Co; 1977.
- 970 63. Tufte ER. *The visual display of quantitative information.* 2nd ed. Cheshire, Conn: Graphics Press; 2001.
- 971 64. Li G-W, Oh E, Weissman JS. The anti-Shine-Dalgarno sequence drives translational pausing and codon choice in  
972 bacteria. *Nature.* 2012;484: 538–541. doi:10.1038/nature10965
- 973 65. Gamble CE, Brule CE, Dean KM, Fields S, Grayhack EJ. Adjacent Codons Act in Concert to Modulate Translation  
974 Efficiency in Yeast. *Cell.* 2016;166: 679–690. doi:10.1016/j.cell.2016.05.070
- 975 66. Chartron JW, Hunt KCL, Frydman J. Cotranslational signal-independent SRP preloading during membrane targeting.  
976 *Nature.* 2016;536: 224–228. doi:10.1038/nature19309
- 977 67. Van Rossum G, Drake Jr FL. *Python reference manual.* Centrum voor Wiskunde en Informatica Amsterdam; 1995.

- 978 68. Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, et al. Biopython: freely available Python tools for  
979 computational molecular biology and bioinformatics. *Bioinformatics*. 2009;25: 1422–1423.  
980 doi:10.1093/bioinformatics/btp163
- 981 69. van der Walt S, Colbert SC, Varoquaux G. The NumPy Array: A Structure for Efficient Numerical Computation.  
982 *Comput Sci Eng*. 2011;13: 22–30. doi:10.1109/MCSE.2011.37
- 983 70. Tange O. GNU Parallel - The Command-Line Power Tool. *Login USENIX Mag*. 2011;36: 42–47.  
984 doi:10.5281/zenodo.16303
- 985 71. R Core Team. R: A Language and Environment for Statistical Computing [Internet]. Vienna, Austria: R Foundation for  
986 Statistical Computing; 2015. Available: <https://www.R-project.org/>
- 987 72. RStudio Team. RStudio: Integrated Development Environment for R. Boston, MA: RStudio, Inc.; 2015.
- 988 73. Bache SM, Wickham H. magrittr: A Forward-Pipe Operator for R [Internet]. 2014. Available: [https://CRAN.R-](https://CRAN.R-project.org/package=magrittr)  
989 [project.org/package=magrittr](https://CRAN.R-project.org/package=magrittr)
- 990 74. Wickham H. The Split-Apply-Combine Strategy for Data Analysis. *J Stat Softw*. 2011;40: 1–29.
- 991 75. Wickham H, Francois R. dplyr: A Grammar of Data Manipulation [Internet]. 2015. Available: [http://CRAN.R-](http://CRAN.R-project.org/package=dplyr)  
992 [project.org/package=dplyr](http://CRAN.R-project.org/package=dplyr)
- 993 76. Aho K. asbio: A Collection of Statistical Tools for Biologists [Internet]. 2015. Available: [http://CRAN.R-](http://CRAN.R-project.org/package=asbio)  
994 [project.org/package=asbio](http://CRAN.R-project.org/package=asbio)
- 995 77. Weinert K. datamart: Unified access to your data sources [Internet]. 2014. Available: [http://CRAN.R-](http://CRAN.R-project.org/package=datamart)  
996 [project.org/package=datamart](http://CRAN.R-project.org/package=datamart)
- 997 78. Wickham H. ggplot2: elegant graphics for data analysis [Internet]. Springer New York; 2009. Available:  
998 <http://had.co.nz/ggplot2/book>
- 999 79. Clarke E, Sherrill-Mix S. ggbeeswarm: Categorical Scatter (Violin Point) Plots [Internet]. 2015. Available:  
1000 <https://github.com/eclarke/ggbeeswarm>
- 1001 80. Schloerke B, Crowley J, Cook D, Briatte F, Marbach M, Thoen E, et al. GGally: Extension to “ggplot2” [Internet]. 2016.  
1002 Available: <https://CRAN.R-project.org/package=GGally>
- 1003 81. Auguie B. gridExtra: Miscellaneous Functions for “Grid” Graphics [Internet]. 2015. Available: [http://CRAN.R-](http://CRAN.R-project.org/package=gridExtra)  
1004 [project.org/package=gridExtra](http://CRAN.R-project.org/package=gridExtra)
- 1005 82. Wilke CO. cowplot: Streamlined Plot Theme and Plot Annotations for “ggplot2” [Internet]. 2015. Available:  
1006 <http://CRAN.R-project.org/package=cowplot>
- 1007 83. Wickham H. scales: Scale Functions for Visualization [Internet]. 2015. Available: [http://CRAN.R-](http://CRAN.R-project.org/package=scales)  
1008 [project.org/package=scales](http://CRAN.R-project.org/package=scales)
- 1009 84. Garnier S. viridis: Default Color Maps from “matplotlib” [Internet]. 2016. Available: [https://CRAN.R-](https://CRAN.R-project.org/package=viridis)  
1010 [project.org/package=viridis](https://CRAN.R-project.org/package=viridis)
- 1011 85. Neuwirth E. RColorBrewer: ColorBrewer Palettes [Internet]. 2014. Available: [http://CRAN.R-](http://CRAN.R-project.org/package=RColorBrewer)  
1012 [project.org/package=RColorBrewer](http://CRAN.R-project.org/package=RColorBrewer)
- 1013 86. Harrower M, Brewer CA. ColorBrewer.org: an online tool for selecting colour schemes for maps. *Cartogr J*. 2003;40:  
1014 27–37.

- 1015 87. Wickham H. multidplyr: Partitioned data frames for “dplyr” [Internet]. Available: <https://github.com/hadley/multidplyr>
- 1016 88. Revolution Analytics, Weston S. foreach: Provides Foreach Looping Construct for R [Internet]. 2015. Available:  
1017 <http://CRAN.R-project.org/package=foreach>
- 1018 89. Revolution Analytics, Weston S. iterators: Provides Iterator Construct for R [Internet]. 2015. Available:  
1019 <https://CRAN.R-project.org/package=iterators>
- 1020 90. Revolution Analytics, Weston S. doMC: Foreach Parallel Adaptor for “parallel” [Internet]. 2015. Available:  
1021 <http://CRAN.R-project.org/package=doMC>
- 1022 91. Revolution Analytics, Weston S. doParallel: Foreach Parallel Adaptor for the “parallel” Package [Internet]. 2015.  
1023 Available: <https://CRAN.R-project.org/package=doParallel>
- 1024 92. Wickham H, Danenberg P, Eugster M. roxygen2: In-Source Documentation for R [Internet]. 2015. Available:  
1025 <https://CRAN.R-project.org/package=roxygen2>
- 1026 93. Zhou J, Rudd KE. EcoGene 3.0. *Nucleic Acids Res.* 2013;41: D613-624. doi:10.1093/nar/gks1235
- 1027 94. Schuler GD, Epstein JA, Ohkawa H, Kans JA. Entrez: molecular biology database and retrieval system. *Methods*  
1028 *Enzymol.* 1996;266: 141–162.
- 1029 95. UniProt Consortium. Reorganizing the protein space at the Universal Protein Resource (UniProt). *Nucleic Acids Res.*  
1030 2012;40: D71-75. doi:10.1093/nar/gkr981
- 1031 96. Cock P. [BioPython] Downloading CDS sequences [Internet]. 2009. Available:  
1032 <http://biopython.org/pipermail/biopython/2009-January/004886.html>
- 1033 97. Tomb JF, White O, Kerlavage AR, Clayton RA, Sutton GG, Fleischmann RD, et al. The complete genome sequence of  
1034 the gastric pathogen *Helicobacter pylori*. *Nature.* 1997;388: 539–547. doi:10.1038/41483
- 1035 98. Lechat P, Hummel L, Rousseau S, Moszer I. GenoList: an integrated environment for comparative analysis of microbial  
1036 genomes. *Nucleic Acids Res.* 2008;36: D469-474. doi:10.1093/nar/gkm1042
- 1037 99. Cole ST, Brosch R, Parkhill J, Garnier T, Churcher C, Harris D, et al. Deciphering the biology of *Mycobacterium*  
1038 *tuberculosis* from the complete genome sequence. *Nature.* 1998;393: 537–544. doi:10.1038/31159
- 1039 100. Dobrovetsky E, Lu ML, Andorn-Broza R, Khutoreskaya G, Bray JE, Savchenko A, et al. High-throughput production  
1040 of prokaryotic membrane proteins. *J Struct Funct Genomics.* 2005;6: 33–50. doi:10.1007/s10969-005-1363-5
- 1041 101. Nelson KE, Clayton RA, Gill SR, Gwinn ML, Dodson RJ, Haft DH, et al. Evidence for lateral gene transfer between  
1042 Archaea and bacteria from genome sequence of *Thermotoga maritima*. *Nature.* 1999;399: 323–329.  
1043 doi:10.1038/20601
- 1044 102. Peden JF. Analysis of codon usage. University of Nottingham. 2000.
- 1045 103. dos Reis M, Wernisch L, Savva R. Unexpected correlations between gene expression and codon usage bias from  
1046 microarray data for the whole *Escherichia coli* K-12 genome. *Nucleic Acids Res.* 2003;31: 6976–6985.
- 1047 104. Zadeh JN, Steenberg CD, Bois JS, Wolfe BR, Pierce MB, Khan AR, et al. NUPACK: Analysis and design of nucleic  
1048 acid systems. *J Comput Chem.* 2011;32: 170–173. doi:10.1002/jcc.21596
- 1049 105. Lorenz R, Bernhart SH, Höner Zu Siederdisen C, Tafer H, Flamm C, Stadler PF, et al. ViennaRNA Package 2.0.  
1050 *Algorithms Mol Biol AMB.* 2011;6: 26. doi:10.1186/1748-7188-6-26

- 1051 106. Coleman JR, Papamichail D, Skiena S, Fitcher B, Wimmer E, Mueller S. Virus attenuation by genome-scale changes  
1052 in codon pair bias. *Science*. 2008;320: 1784–1787. doi:10.1126/science.1155761
- 1053 107. Yang ZR, Thomson R, McNeil P, Esnouf RM. RONN: the bio-basis function neural network technique applied to the  
1054 detection of natively disordered regions in proteins. *Bioinformatics*. 2005;21: 3369–3376.  
1055 doi:10.1093/bioinformatics/bti534
- 1056 108. Kuhn M. Building predictive models in R using the caret package. *J Stat Soft*. 2008;
- 1057 109. Weihs C, Ligges U, Luebke K, Raabe N. klaR Analyzing German Business Cycles. In: Baier D, Decker R, Schmidt-  
1058 Thieme L, editors. *Data Analysis and Decision Support*. Berlin/Heidelberg: Springer-Verlag; 2005. pp. 335–343.  
1059 Available: [http://link.springer.com/10.1007/3-540-28397-8\\_36](http://link.springer.com/10.1007/3-540-28397-8_36)
- 1060 110. Kendall MG. A New Measure of Rank Correlation. *Biometrika*. 1938;30: 81. doi:10.2307/2332226
- 1061 111. Joachims T. Optimizing search engines using clickthrough data. ACM Press; 2002. p. 133. doi:10.1145/775047.775067
- 1062 112. Robin X, Turck N, Hainard A, Tiberti N, Lisacek F, Sanchez J-C, et al. pROC: an open-source package for R and S+ to  
1063 analyze and compare ROC curves. *BMC Bioinformatics*. 2011;12: 77. doi:10.1186/1471-2105-12-77
- 1064 113. DeLong ER, DeLong DM, Clarke-Pearson DL. Comparing the areas under two or more correlated receiver operating  
1065 characteristic curves: a nonparametric approach. *Biometrics*. 1988;44: 837–845.
- 1066 114. Freedman D, Diaconis P. On the histogram as a density estimator:L 2 theory. *Z F  $\diamond$  r Wahrscheinlichkeitstheorie*  
1067 *Verwandte Geb*. 1981;57: 453–476. doi:10.1007/BF01025868
- 1068 115. Cunningham F, Amode MR, Barrell D, Beal K, Billis K, Brent S, et al. Ensembl 2015. *Nucleic Acids Res*. 2015;43:  
1069 D662-669. doi:10.1093/nar/gku1010
- 1070 116. Smedley D, Haider S, Durinck S, Pandini L, Provero P, Allen J, et al. The BioMart community portal: an innovative  
1071 alternative to large, centralized data repositories. *Nucleic Acids Res*. 2015;43: W589-598. doi:10.1093/nar/gkv350
- 1072 117. Harris TW, Baran J, Bieri T, Cabunoc A, Chan J, Chen WJ, et al. WormBase 2014: new views of curated biology.  
1073 *Nucleic Acids Res*. 2014;42: D789-793. doi:10.1093/nar/gkt1063
- 1074 118. Engel SR, Dietrich FS, Fisk DG, Binkley G, Balakrishnan R, Costanzo MC, et al. The reference genome sequence of  
1075 *Saccharomyces cerevisiae*: then and now. *G3 Bethesda Md*. 2014;4: 389–398. doi:10.1534/g3.113.008995
- 1076 119. De Schutter K, Lin Y-C, Tiels P, Van Hecke A, Glinka S, Weber-Lehmann J, et al. Genome sequence of the  
1077 recombinant protein production host *Pichia pastoris*. *Nat Biotechnol*. 2009;27: 561–566. doi:10.1038/nbt.1544
- 1078 120. Nordberg H, Cantor M, Dusheyko S, Hua S, Poliakov A, Shabalov I, et al. The genome portal of the Department of  
1079 Energy Joint Genome Institute: 2014 updates. *Nucleic Acids Res*. 2014;42: D26-31. doi:10.1093/nar/gkt1069
- 1080 121. Markowitz VM, Chen I-MA, Palaniappan K, Chu K, Szeto E, Pillay M, et al. IMG 4 version of the integrated  
1081 microbial genomes comparative analysis system. *Nucleic Acids Res*. 2014;42: D560-567. doi:10.1093/nar/gkt963
- 1082 122. Uchiyama I, Mihara M, Nishide H, Chiba H. MGD update 2015: microbial genome database for flexible ortholog  
1083 analysis utilizing a diverse set of genomic data. *Nucleic Acids Res*. 2015;43: D270-276. doi:10.1093/nar/gku1152
- 1084 123. Loader C. locfit: Local Regression, Likelihood and Density Estimation. [Internet]. 2013. Available: [https://CRAN.R-](https://CRAN.R-project.org/package=locfit)  
1085 [project.org/package=locfit](https://CRAN.R-project.org/package=locfit)
- 1086 124. Schindelin J, Arganda-Carreras I, Frise E, Kaynig V, Longair M, Pietzsch T, et al. Fiji: an open-source platform for  
1087 biological-image analysis. *Nat Methods*. 2012;9: 676–682. doi:10.1038/nmeth.2019



- 1088 125. Schneider CA, Rasband WS, Eliceiri KW. NIH Image to ImageJ: 25 years of image analysis. *Nat Methods*. 2012;9:  
1089 671–675.
- 1090 126. Xie Y. knitr: A Comprehensive Tool for Reproducible Research in R. In: Stodden V, Leisch F, Peng RD, editors.  
1091 Implementing Reproducible Computational Research. Chapman and Hall/CRC; 2014. Available:  
1092 <http://www.crcpress.com/product/isbn/9781466561595>
- 1093 127. Towns J, Cockerill T, Dahan M, Foster I, Gauthier K, Grimshaw A, et al. XSEDE: Accelerating Scientific Discovery.  
1094 *Comput Sci Eng*. 2014;16: 62–74. doi:10.1109/MCSE.2014.80

1095

## 1096 **Supporting Information**

1097 **S1 Fig. Additional small-scale predictions and outcomes.** (A) Experimental expression of 116 *H.*  
1098 *pylori* membrane proteins in *E. coli* in at most 3 vectors (238 trials) scored as either a 1, 2, or 3 from the  
1099 outcome of a dot blot as well as Coomassie Staining of an SDS-PAGE gel for two of the vectors. To  
1100 compare the three vectors with a single set of scores, the two scores were averaged to give a single  
1101 number for a condition making them comparable to the third vector while yielding 2 additional  
1102 thresholds (1.5 and 2.5) and the 6 total levels shown. (B) The Receiver Operating Characteristic (ROC)  
1103 with each cutoff is plotted, where a higher cutoff is represented by a deeper red, followed by the Area  
1104 Under the Curves (directly below) in colors that correspond to the respective curve. (C) Expression of  
1105 77 *T. maritima* membrane proteins in *E. coli* noted as purified (5), not purified but expressed (14), or  
1106 neither. (D) ROC curve for each threshold. (E) Expression of 37 microbial secondary transporters in 4  
1107 IPTG-inducible vectors (144 trials) in *E. coli* quantified as 10 ng/mL (pink) or 100 ng/mL (red) via dot  
1108 blot. (F) ROC curve for each threshold.

1109 **S2 Fig. Complete set of forward predictions.** (A) Extended from Fig 5C, the full complement of score  
1110 distributions calculated by genome is plotted and arranged to accentuate similar features by physiology,  
1111 *e.g.* growth condition, or scientific interest, *e.g.* pathogenic. Raw scores along with sequence identifiers  
1112 are available in the S4 Table. (B) Histograms of representative sequence features between the training  
1113 data set (green), thermophiles (orange), and *P. falciparum* (purple). Values for sequence parameter  
1114 overlap coefficients derived from kernel density estimates (Methods 7) versus the *E. coli* training data  
1115 are included. See S1 Table for parameter descriptions.

1116 **S3 Fig. Complete set of feature correlations and their individual contributions to the model.**  
1117 Features are ordered first by category (as in Fig 5) and then by weight (grey bars). Labels are green for  
1118 protein-sequence derived and brown for nucleotide-sequence derived features. Pearson correlation  
1119 coefficient between each pair of features across the NYCOMPS dataset is plotted (right). See S1 Table  
1120 for a detailed description of each feature. Feature categories are overlaid as square boxes and indicated  
1121 by black bars on the top, left, and right of the correlation matrix.

1122 **S4 Fig. Feature contributions to the model across datasets used for training and validation.** (A)  
1123 Total weight for each category is represented as a bar. The contribution of each feature to the category is  
1124 shown by partitioning the bar. The red dot indicates the total sum of weights within the category. (B)  
1125 Pearson correlation coefficients between feature categories are shown. Feature labels are green for  
1126 protein-sequence derived and brown for nucleotide-sequence derived. (C) Feature category dependence  
1127 within the training set is shown by Spearman's  $\rho$  and 95% CI between the normalized outcomes versus

1128 the feature subset. **(D)** Considering the NYCOMPS data set (as in Fig 2), the Area Under the Curve  
1129 (AUC) of a Receiver Operating Characteristic and 95% confidence interval when predicting solely by  
1130 features from the specified category against the NYCOMPS dataset. Red, using positive only as the cut-  
1131 off for individual genes (Fig 2C); grey, using positive outcomes within each plasmid and solubilization  
1132 condition (as in Fig 2E).

1133 **S1 Table. Sequence parameter weights and descriptions.** Weights are presented after normalizing to  
1134 the mean value for clarity. Features that were calculated but removed in pre-processing are noted  
1135 (Methods 3).

1136 **S2 Table. AUC values for the NYCOMPS dataset.** AUC values and 95% confidence intervals are  
1137 presented in summary, by expression condition, and by predicted C-terminal localization as well as for  
1138 IMProve scores calculated without the most computationally expensive RNA secondary structure  
1139 calculation (as in Fig 5).

1140 **S3 Table. Predictive performances of the model across protein families.** The proteins and  
1141 performances are with respect to those tested by NYCOMPS as summarized in Fig 5. This data is  
1142 available in an interactive format at [clemonslab.caltech.edu](http://clemonslab.caltech.edu).

1143 **S4 Table. Full list of predicted membrane proteins.** This includes corresponding identifiers,  
1144 descriptions, Pfam families, coding sequences, and IMProve scores. This data is available in an  
1145 interactive format at [clemonslab.caltech.edu](http://clemonslab.caltech.edu).

Figure 1

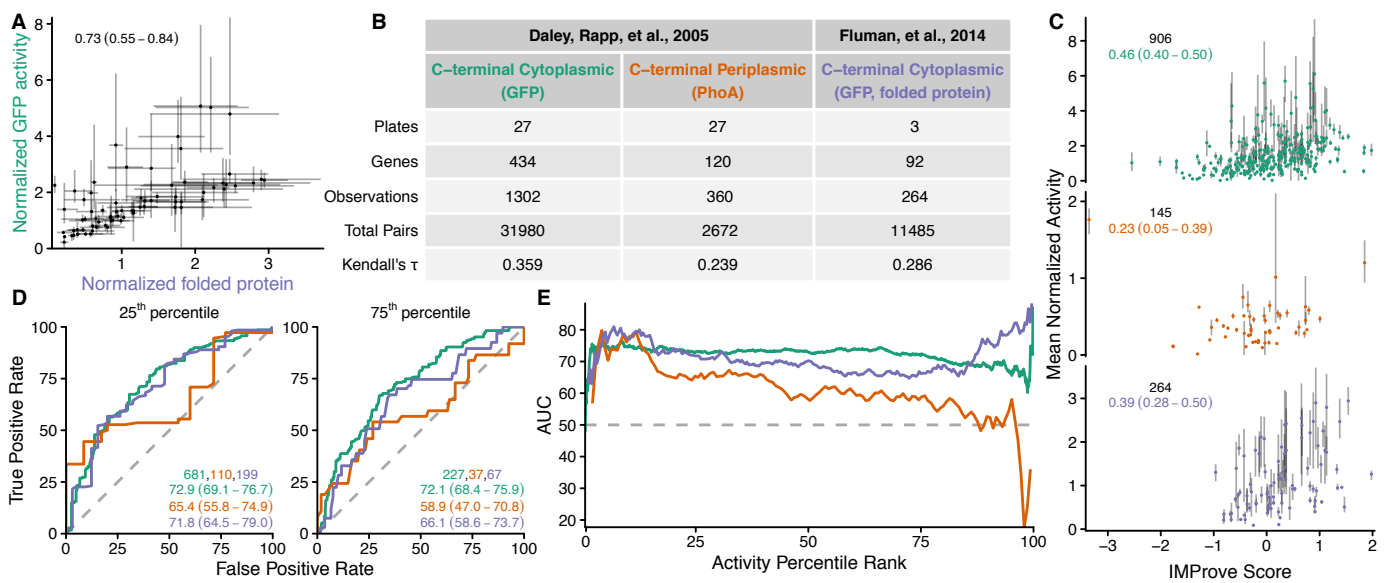


Figure 2

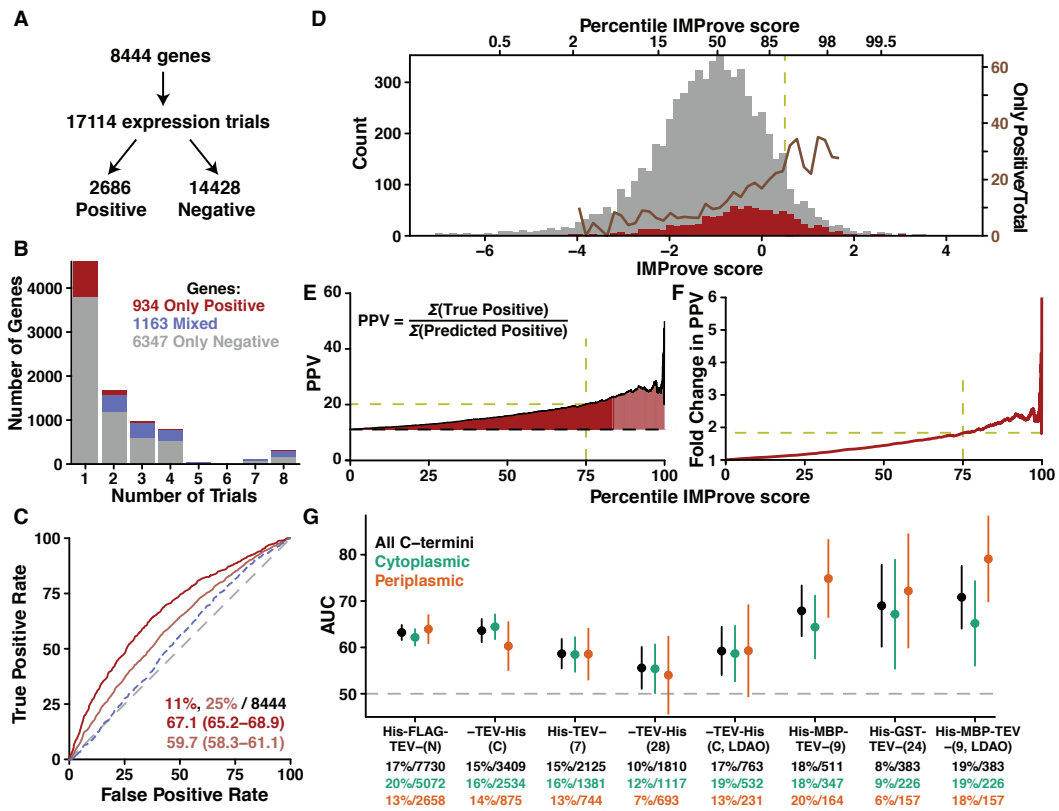


Figure 3

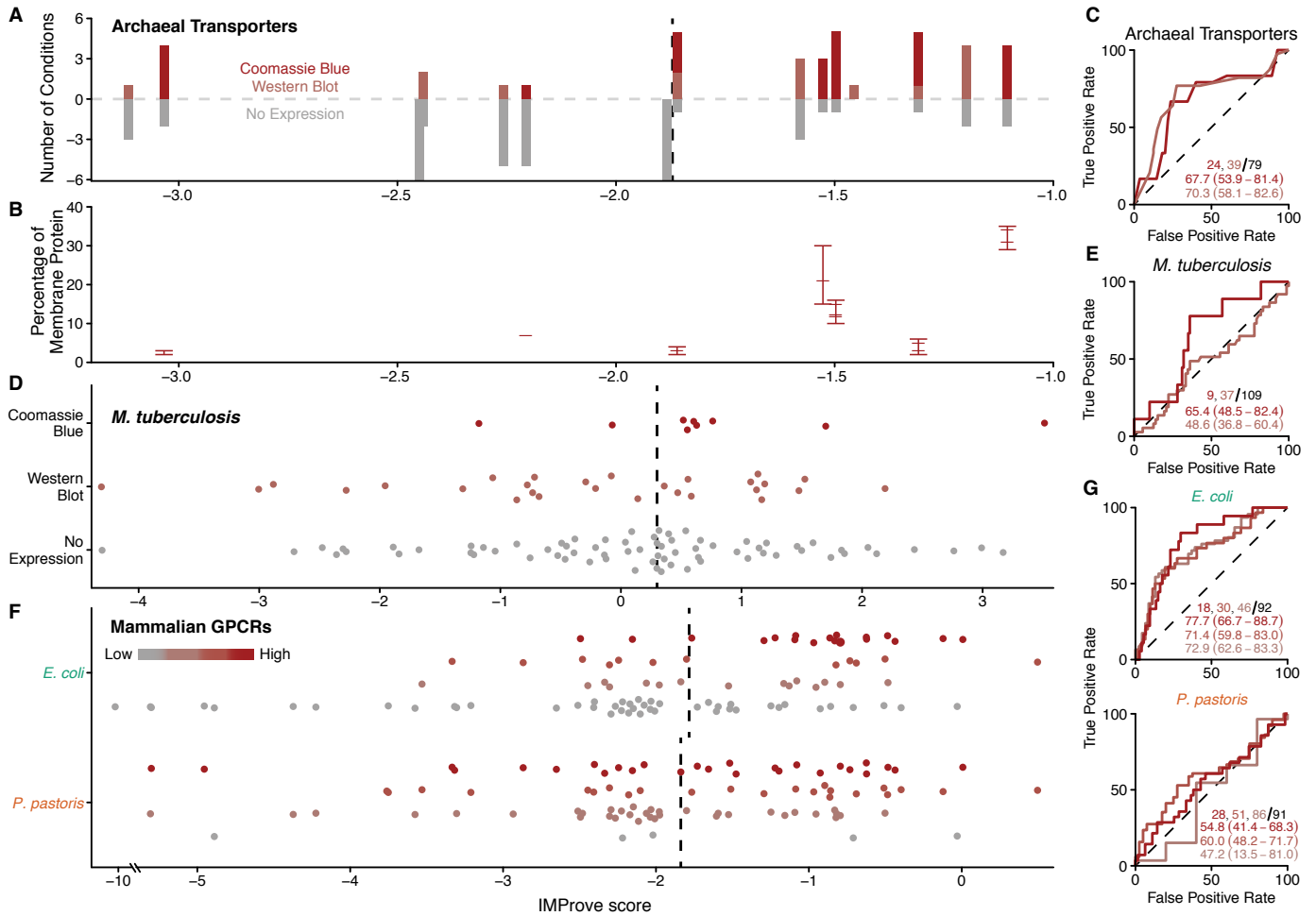




Figure 4

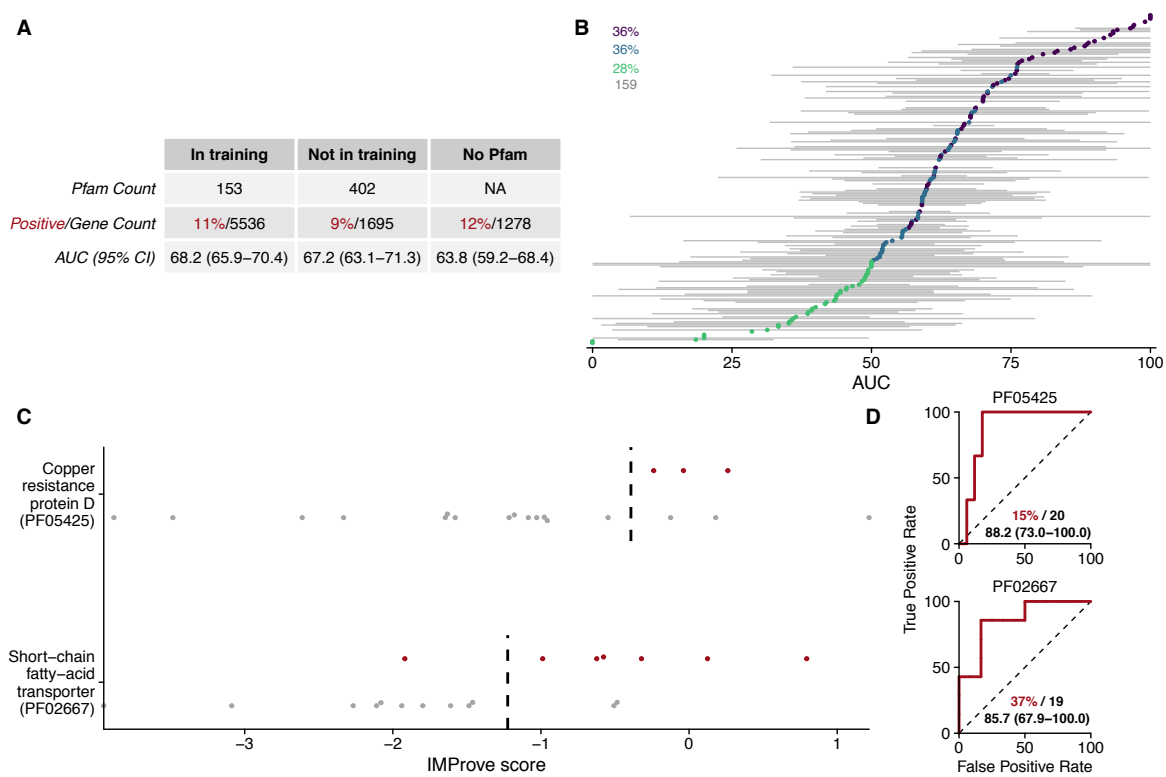


Figure 5

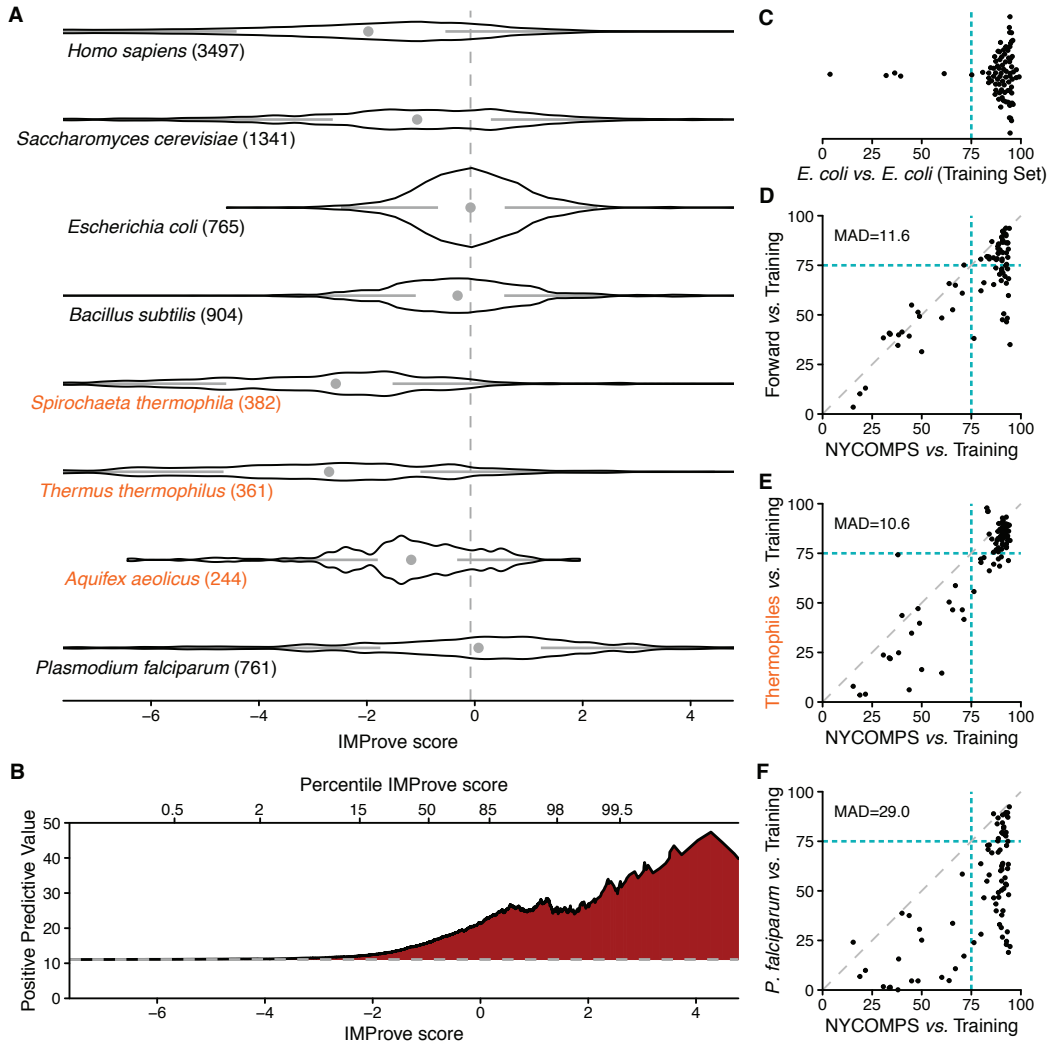


Figure 6

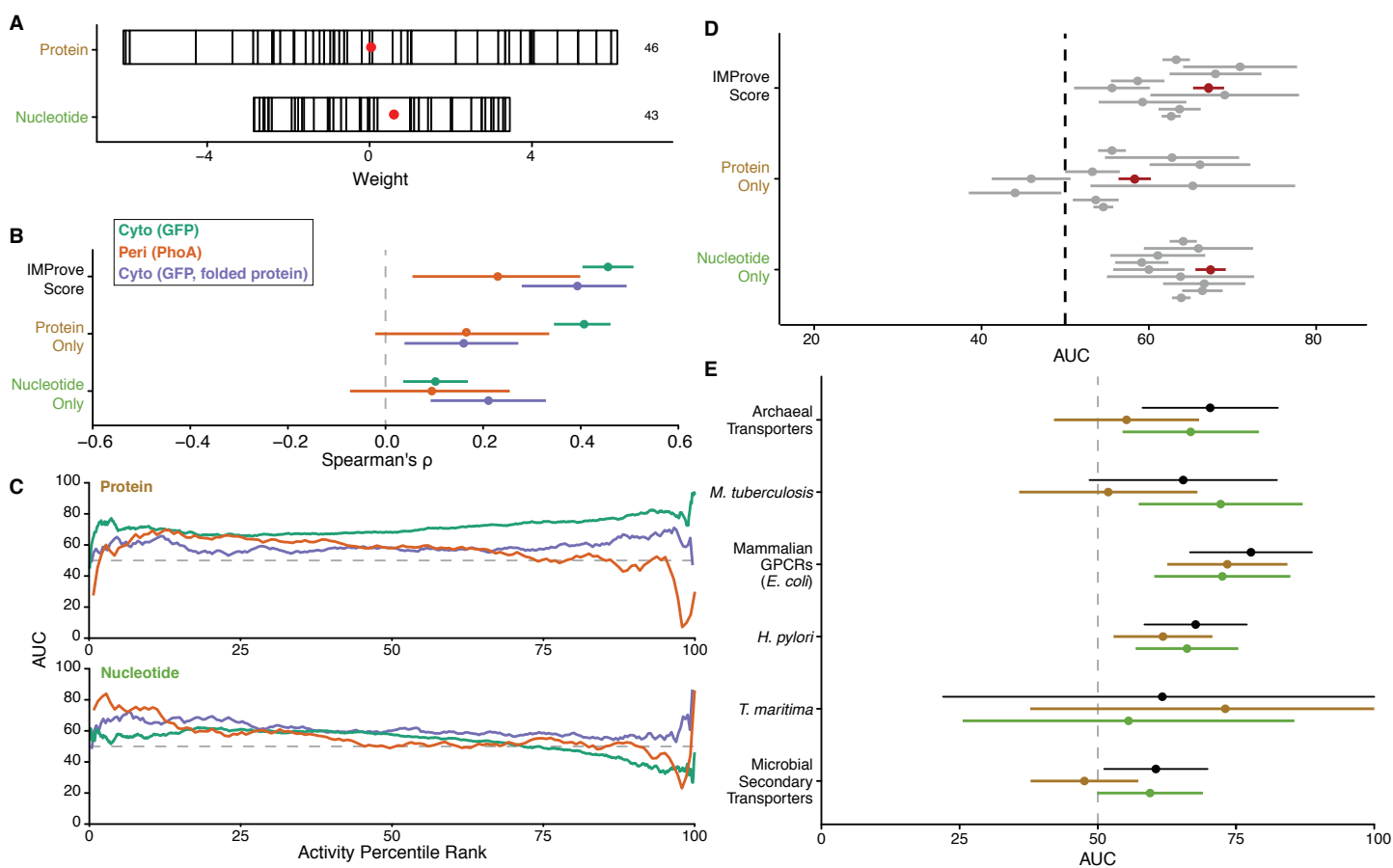
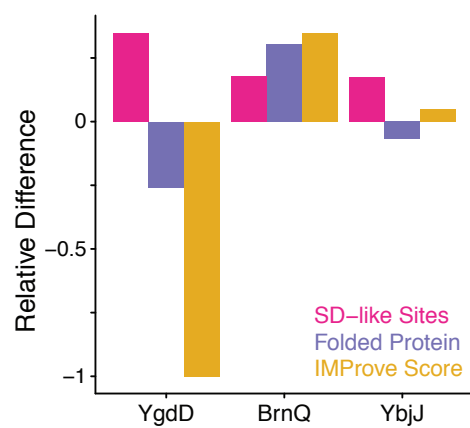
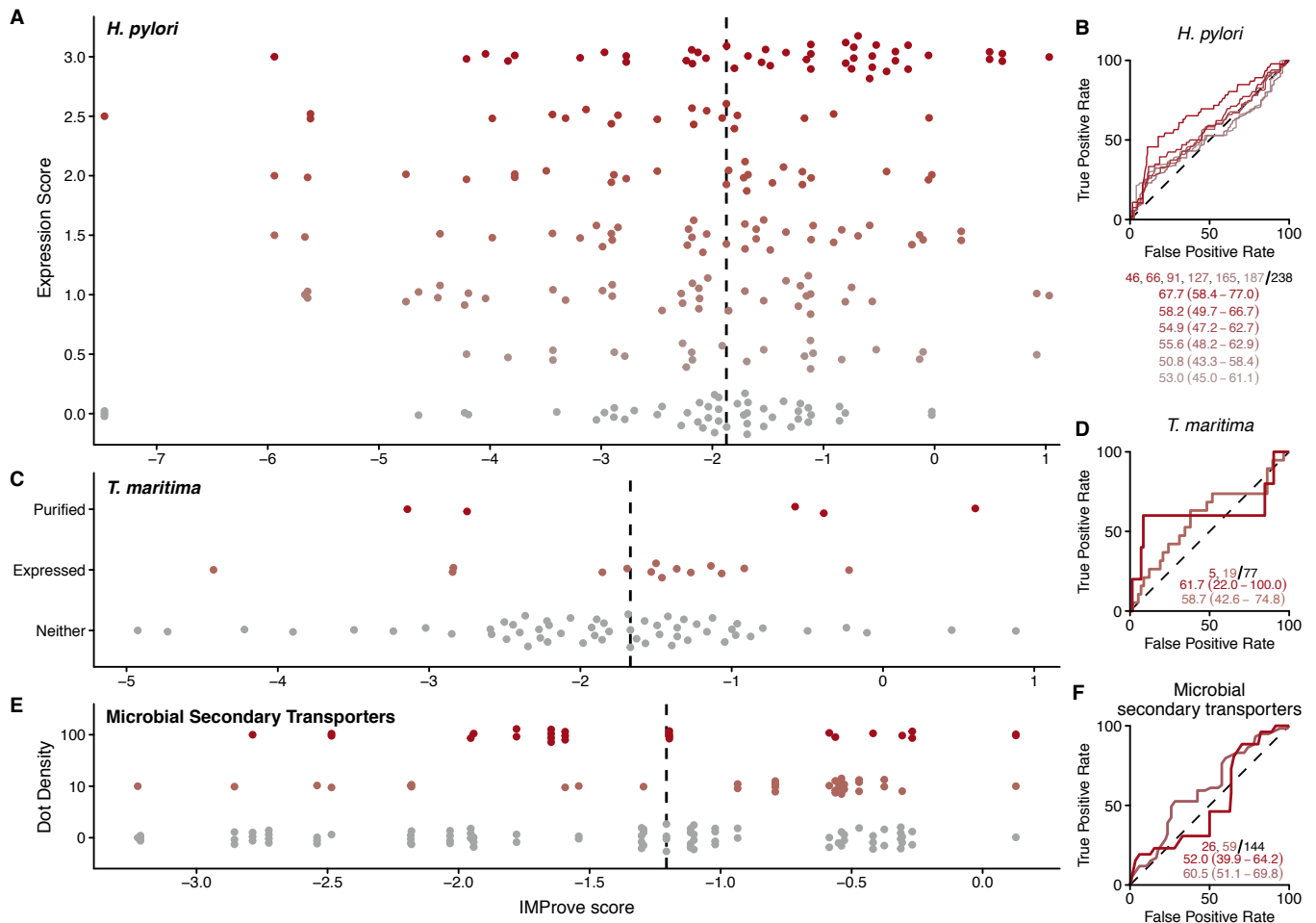


Figure 7

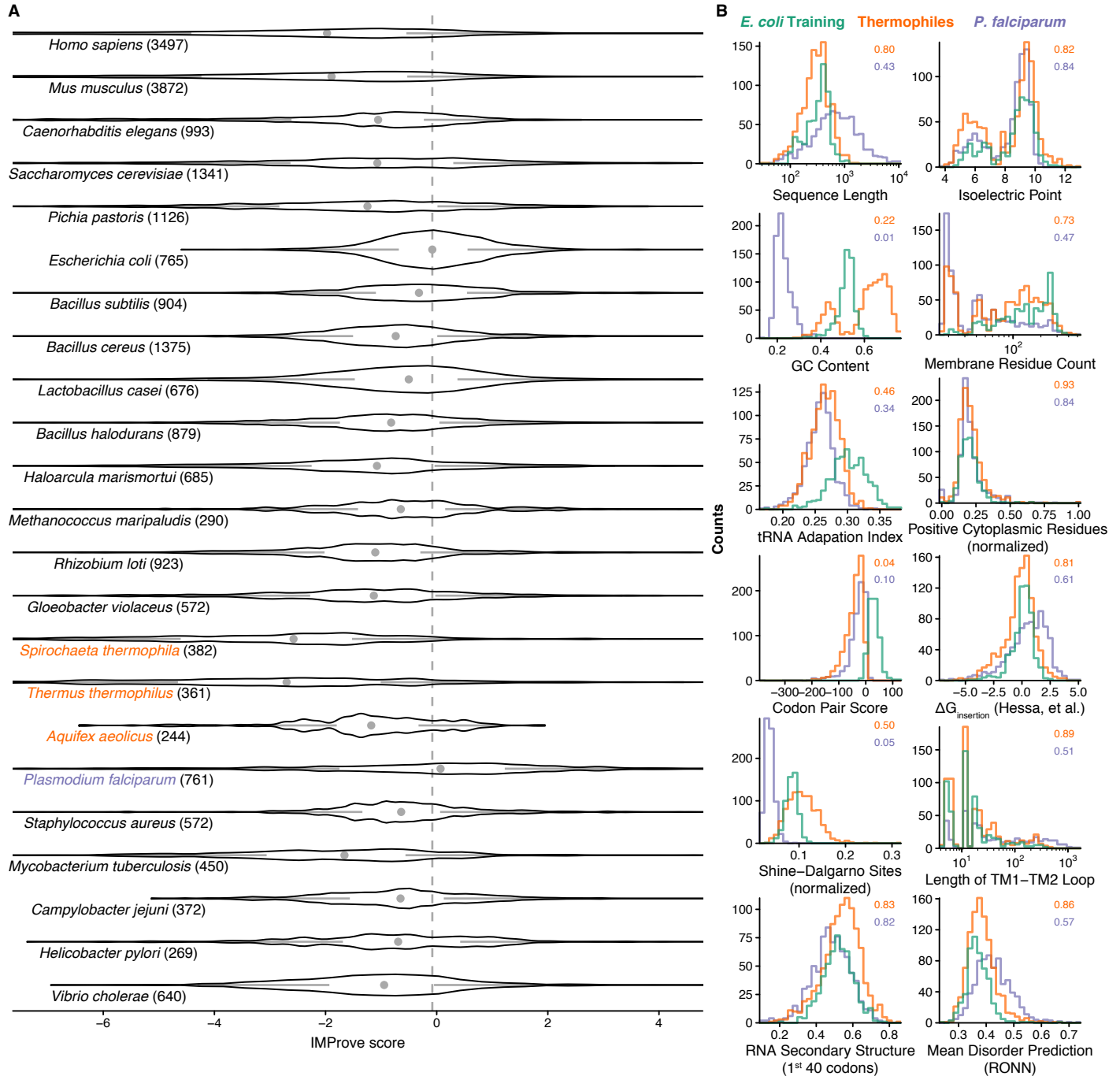


## S1 Figure

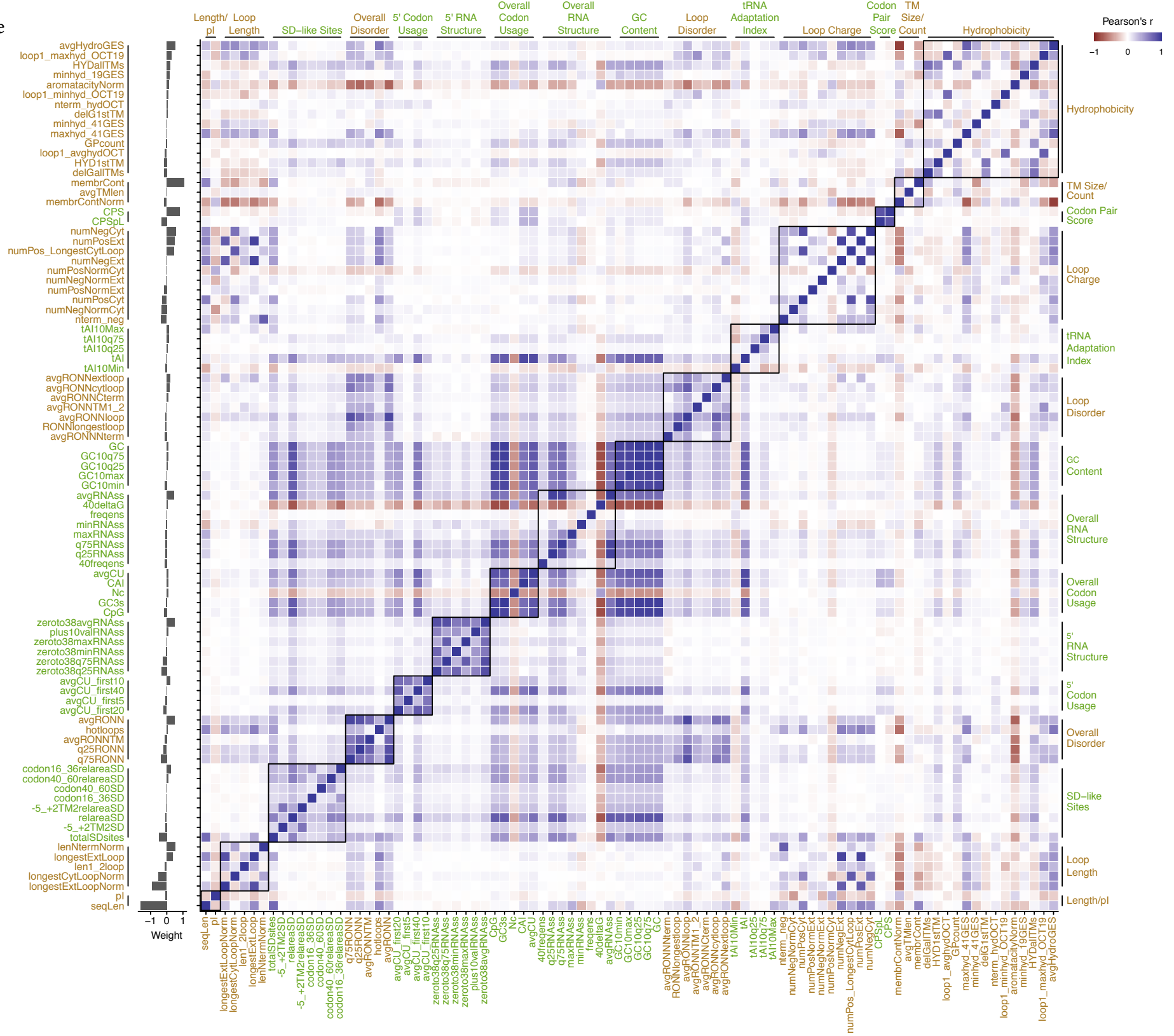




S2 Figure



# S3 Figure



## S4 Figure

