



23 and activity. For example, solubility of a tyrosine ammonium lyase was more than dou-  
24 bled by adding two tags to its N- and C-terminus. Its protein activity was also increased  
25 nearly 3.5 fold by adding the tags. Additional experiments also supported that the de-  
26 signed tags were effective for improving activity of multiple proteins and are better than  
27 previously reported tags. The presented optimization methodology thus provides a val-  
28 uable tool for understanding the correlation between amino acid sequence and protein  
29 solubility and for engineering protein biocatalysts.

30 **Contact:** kang.zhou@nus.edu.sg, chewxia@nus.edu.sg

31

## 32 **Introduction**

33 The exploration of expressing recombinant proteins started in 1976, when human pep-  
34 tide hormone Somatostatin was produced in *Escherichia coli*<sup>1</sup>. As the most commonly  
35 used expression host, *E. coli* was investigated intensively to improve the expression  
36 and activity of recombinant proteins<sup>2,3,4</sup>. Various experimental strategies, such as using  
37 protein fusion partners, co-expressing chaperones, choosing suitable promoters, opti-  
38 mizing codon usage, changing culture conditions, or using directed evolution<sup>5, 6, 7, 8, 9,</sup>  
39 <sup>10</sup>, were used to improve protein expression. For example, the expression of human  
40 recombinant enzyme N-acetylgalactosamine-6-sulfatase (rhGALNS) in *E. coli* was un-  
41 desirable due to protein aggregation. Several methods including the use of physiologi-  
42 cally-regulated promoters, overexpression of native chaperones and applying osmotic  
43 shock were investigated to improve the production and activity of rhGALNS<sup>10</sup>. Protein  
44 activity, a phenotype representing the catalytic ability of a protein if it is an enzyme, is  
45 partly determined by its genotype (sequence of its coding gene). Directed evolution can  
46 effectively improve protein activity through changing the associated genotype, but this

47 approach is resources-intensive. In the process of improving protein activity via di-  
48 rected evolution, mutagenesis is performed to change gene sequence and the mutated  
49 genes are inserted into plasmid used for transformation of a microbe, usually *E. coli*.  
50 Additional techniques are employed further to screen a large number of transformed  
51 cells for those that have higher protein activity. Since most of the protein directed evo-  
52 lution studies were only interested in the mutants with the highest activity, they did not  
53 reveal the genotype of most proteins that had lower activity. This fact has caused the  
54 challenge that almost no suitable database of protein activity is available for training  
55 computational models that can predict protein activity from protein sequence. Such  
56 models would greatly assist protein engineering by evaluating protein sequences *in sil-*  
57 *ico*. A suitable dataset for training the model should contain both protein activity data  
58 and the associated sequence data, and should be large enough (>1,000 entries).

59 Protein activity data cannot be easily pooled together for model training if they are  
60 related to enzymes that catalyze different chemistries, which is another reason why it  
61 is difficult to generate the aforementioned datasets. The data of protein solubility from  
62 most types of proteins, however, can be compiled into one dataset, because protein sol-  
63 ubility is a basic protein property. In this study and the relevant literature, protein solu-  
64 bility is defined as the percentage of a protein's soluble fraction<sup>11</sup>. It is a metric that is  
65 often used to assess the folding quality of a protein, under the assumption that incor-  
66 rectly folded proteins form aggregates and are insoluble. Protein activity is thus corre-  
67 lated with protein solubility to some extent, because protein solubility may indicate the  
68 quality of protein folding which influences protein 3D structure and activity, i.e. pro-  
69 teins with higher solubility likely exhibit higher activity<sup>12</sup>. Improving the solubility of  
70 some recombinant proteins can enhance their production effectively<sup>13</sup>. Thus, protein  
71 solubility may be used as a proxy for protein activity to develop predictive models that

72 use protein sequence as input. With such a model, it would be possible to optimize the  
73 protein sequence of a protein *in silico* for improving its solubility and activity. For ex-  
74 ample, a Monte Carlo optimization method can be used as the procedures demonstrated  
75 in Figure 1: (1) a random change is introduced to the protein sequence, (2) the new  
76 protein sequence is evaluated by the model, and (3) if the predicted solubility is lower  
77 than that of the parental sequence, the change would be rejected, otherwise it would be  
78 accepted and used to initiate the subsequent iteration. This *in silico* optimization pro-  
79 cess may identify promising protein sequences to improve the success rate of the time-  
80 consuming and labor-intensive experiments. If the protein activity heavily depends on  
81 its solubility, the experiment would identify new protein that has higher solubility and  
82 activity.

83 Machine learning has gained increasing attention recently in various fields, such as in-  
84 ternet commerce, autonomous vehicles, and image recognition<sup>14, 15, 16, 17, 18, 19, 20, 21, 22</sup>.  
85 Until now, a large number of machine learning methods have been explored to predict  
86 protein solubility from amino acid sequence<sup>6, 11, 23, 24, 25</sup>. Among the previous studies,  
87 we developed regression models that can predict protein solubility in the continuous  
88 values<sup>26</sup>. Classification models which only label a protein as soluble or insoluble were  
89 developed in other studies but cannot be used in the *in silico* optimization, because it  
90 would mistakenly reject most changes that can result in a small but important increase  
91 in the protein solubility. So far, very few studies performed experimental validation of  
92 their solubility-prediction models and no study used such models to improve protein  
93 properties through the *in silico* optimization of protein sequence.

94 In our present study, based on a regression model that can predict protein solubility  
95 from protein sequence<sup>26</sup>, we developed optimization algorithms to increase predicted  
96 solubility under constraints that have been set after considering experimental feasibility

97 and impact on protein function. The performance of the optimization process for im-  
98 proving protein solubility was validated successfully by experimentally measuring sol-  
99 ubility. We found that adding short peptide rich in negatively charged amino acids was  
100 effective in improving solubility of many proteins. More importantly, we also verified  
101 that activity of some proteins was indeed substantially improved when their solubility  
102 was increased. Our study provides a generally effective approach to enhance protein  
103 solubility and activity.

104

## 105 **Results**

### 106 **Design the optimization methodology**

107 In order to improve protein solubility by *in silico* mutagenesis, we need to solve several  
108 questions regarding how to change the protein sequence. One can change a protein se-  
109 quence by adding amino acids to the sequence (addition), replacing amino acids in the  
110 sequence (mutation) and/or removing amino acids from the sequence (deletion). The  
111 protein functions may be frequently abolished by mutation and deletion as the original  
112 protein structure and active sites may be changed. To avoid such detrimental change to  
113 the original function of the protein, addition was used in our study to change protein  
114 sequence for improving protein solubility. The subsequent decision to make is how  
115 many amino acids should be added. Adding too many amino acids would make exper-  
116 imental validation to be more expensive and may also negatively affect the protein  
117 function. Adding too few amino acids may not be able to improve protein solubility  
118 substantially. We decided to evaluate adding 20 or 30 amino acids because adding more  
119 than 30 amino acids to a protein by using synthetic oligonucleotides was experimentally  
120 difficult.

121 To optimize the sequence of the amino acids to be added, we designed an algorithm  
122 based on the support-vector machine (SVM) prediction model we previously devel-  
123 oped<sup>26</sup>. The independent variable in the optimization function is the amino acid com-  
124 position of the short peptide to be added, expressed as number of each amino acid in a  
125 vector (Figure 1). The SVM model we developed only accepted amino acid composi-  
126 tion of a protein as input, so we did not consider the full sequence information during  
127 the optimization. Then the amino acid composition of a model protein with the added  
128 amino acids was calculated and used as input for the SVM model. We used the genetic  
129 algorithm (GA) which is a widely used algorithm for solving constrained optimization  
130 problems. The objective function of GA outputs the predicted protein solubility by us-  
131 ing the SVM model in the format of continuous values between 0-1. The sum of the  
132 number of amino acids added was set as 20 or 30 and the searching range for the number  
133 of each amino acid added was from 0 to 20 or 30.

134

### 135 **Optimize protein sequence *in silico* for improving protein solubility**

136 After designing this optimization algorithm, ten proteins with low solubility (0.1) in the  
137 eSol database (we had used the same database to train our machine learning model)  
138 were selected as model proteins to test the algorithm (information of these proteins is  
139 provided in Supplementary Table S2). The predicted solubility of all the ten proteins  
140 was improved after adding 30 amino acids as peptide tags (Supplementary Figure S2).  
141 One protein's solubility (name: agaW, N-acetylgalactosamine-specific enzyme IIC  
142 component of PTS) was improved to 0.9951 from 0.1 after adding the designed short  
143 peptide tags. When we allowed adding only 20 instead of 30 amino acids, the improve-  
144 ment of predicted solubility slightly decreased (Supplementary Figure S2). Since it is

145 easier and cheaper to add 20 amino acids in experiments than 30, we adopted adding  
146 20 amino acids as the constraint in the rest of this study.

147 To make this study more relevant to the imperative applications of recombinant en-  
148 zymes, we selected six proteins which were important in engineering metabolic path-  
149 way of *E. coli* to produce valuable metabolites (information of these proteins is pro-  
150 vided in caption of Figure 2). These proteins' predicted solubility was lower than 0.6.  
151 Adding 20 amino acids also substantially improved the predicted solubility of all the  
152 six proteins (Figure 2). Three proteins (tal, dxs and valC) were chosen to experimentally  
153 validate the optimization results since their original predicted solubility was low and  
154 the predicted solubility was substantially improved through the optimization.

155 We also included agaw in the test because of the large improvement we observed in the  
156 *in silico* optimization. The number of the amino acids to be added was allowed to be  
157 decimal during the optimization and was rounded for experimental validation. The pre-  
158 dicted solubility after rounding the number of the amino acids added was very similar  
159 to that before rounding for all the tested proteins (Supplementary Table S6). To generate  
160 sequence of the two tags to be added to a protein from the number of amino acids we  
161 minimized the occurrence of amino acid repeats, which reduced the difficulty in syn-  
162 thesizing the DNA. The sequence of the tags for those four proteins is listed in the  
163 Supplementary Table S7.

164

### 165 **Experimental validation of the optimized protein sequence**

166 We constructed expression vectors to express the four proteins with and without the  
167 optimized tags. Among them, protein agaw cannot be expressed (as determined by us-  
168 ing SDS-PAGE) with and without the tags, which may be caused by the unstable protein

169 structure or unsuitable experimental conditions. Protein valC can be expressed only  
170 without the peptide tags which may have impaired the protein stability. Protein tal and  
171 dxs were expressed with and without the tags (Figure 3). Protein solubility of tal and  
172 dxs was increased by 118% and 16% respectively by adding the tags.

173 By observing the amino acids added to dxs and tal (Figure 3b and Supplementary table  
174 5), it can be found that their peptide tags were dominated by aspartic acid (D) and glu-  
175 tamic acid (E). Aspartic acid and glutamic acid are the two negatively charged amino  
176 acids among the 20 amino acids. Adding them may introduce repulsive electrostatic  
177 interactions between protein molecules to prevent aggregation and to provide sufficient  
178 time for correct folding of proteins<sup>27</sup>. The similarity of the peptide tags inspired us to  
179 test whether one tag designed for one protein can be used to improve solubility of an-  
180 other protein. We found that the tags optimized for improving solubility of tal could  
181 also increase both predicted and measured solubility of dxs, and vice versa (Figure 4a).  
182 Another protein (name: ada, aldehyde dehydrogenase) used in a project of our labora-  
183 tory was also tested with the tag designed for tal and its predicted and measured solu-  
184 bility were also enhanced (Figure 4a). The results of switching tags suggested that the  
185 tags we designed may be generally effective in improving protein solubility.

186

### 187 **Protein activity also improved by the optimization**

188 The ultimate goal of this project was to improve activity of enzymes and their solubility  
189 was used as proxy because of the aforementioned reasons. Following the success of  
190 improving protein solubility, we measured activity of tal with and without the tags.  
191 Protein tal is tyrosine ammonia lyase which can deaminate tyrosine to produce couma-



192 ric acid (Figure 4c). It is very useful in producing flavonoids by using engineered mi-  
193 crobes<sup>28, 29</sup>. Tal activity was increased by 269% by adding the tags we designed for it  
194 (Figure 4d, based on 12 h reaction). The extent of the increase in activity was even  
195 larger than that in solubility, suggesting that adding the tags may also increase the ex-  
196 pression level and/or specific activity of tal. This result proved that our optimization  
197 scheme for protein solubility was also effective for improving protein activity and using  
198 protein solubility as a proxy to increase protein activity was reasonable.

199

### 200 **Tags designed under more constrained conditions**

201 Among the four proteins selected for experimental validation, the protein valC (valen-  
202 cene synthase) cannot be synthesized only after the tags were added. This may be  
203 caused by the fact that the stability of protein valC was damaged after adding the tags.  
204 Our prediction model and optimization algorithm only took the protein solubility into  
205 account. However, other properties of the protein may be changed during the addition  
206 of highly charged tags, such as the protein stability. Therefore, we explored whether the  
207 peptide tags including mainly aspartic acid and glutamic acid can be replaced by tags  
208 that contain less charged amino acids to improve protein solubility.

209 The constrained condition that the number of aspartic acid and glutamic acid cannot be  
210 more than a threshold was therefore set in the optimization algorithm. The threshold  
211 was from 0 to 10 with step size of 1 for aspartic acid and glutamic acid respectively  
212 (Supplementary Table S8). When the limitation of addition number for aspartic acid  
213 and glutamic acid was reduced gradually from 10, the predicted solubility was decreas-  
214 ing but the change was small. With the decrease in the number of aspartic acid and  
215 glutamic acid, the number of lysine (K) increased substantially. Other amino acids only

216 had a relatively small increase in the optimization solutions. When the constrained con-  
217 dition was very strict, for example, no aspartic acid and glutamic acid were allowed,  
218 the amino acids introduced were mostly alanine (A).

219 Another constrained condition was explored which limited the net charge of the peptide  
220 tags. In this case, the upper bound for the absolute value of the net charge of the tag  
221 was set as 5, 4, 3, 2, 1, and 0, respectively (Supplementary Table S9) and it could be  
222 observed that the number of alanine increased most substantially with the decrease of  
223 net charge, which was consistent with the results obtained under the other constraint  
224 and supported that introducing alanine may be beneficial for the dissolution of protein  
225 or the optimization failed to find a feasible solution under such stringent constraints.

226 This hypothesis was tested by doing experiments. The tags with net charge 1, 3, and 5  
227 (Supplementary table S9) were used with protein valC. These new tags did not abolish  
228 protein expression, confirming the hypothesis that excessive amount of aspartic acid  
229 and/or glutamic acid may destabilize certain proteins. However, the solubility of protein  
230 valC was not improved by the tags (Supplementary Figure S3). Protein valC may have  
231 strong affinity to cellular membranes and thus cannot be solubilized by the designed  
232 tags.

233

#### 234 **Comparison with previous studies**

235 To improve protein solubility, some trial-and-error procedures were developed by in-  
236 troducing small polyionic tags<sup>30, 31, 32</sup>. Small peptide tags have been used as solubility-  
237 enhancing tags for a long time because they are short and less likely to interfere with  
238 protein structure<sup>30</sup>. One study indicated non-polar surface and positively-charged  
239 patches contributed to the separation of the soluble and insoluble proteins<sup>31</sup>. It was

240 demonstrated that a concentration of positive charge may tend towards lower folded  
241 state stability through unfavourable charge interactions and result in insolubility. In ad-  
242 dition, a negatively charged fusion tag, NT11, was also developed to enhance protein  
243 expression in *E. coli*<sup>32</sup>. However, these previous studies explored tags by trial and error  
244 and cannot provide a generally useful quantitative model which can forecast perfor-  
245 mance of tags with proteins which have not been tested. Among the diverse solubility-  
246 enhancing tags that have been tested, the ones that are rich in aspartic acid and glutamic  
247 acid were also studied before<sup>27</sup>.

248 To find out if the tags we obtained from our optimization were more effective than these  
249 published ones, we compared them by using our predictive model and by conducting  
250 experiments. We used tal as the model protein here, because its solubility was experi-  
251 mentally confirmed to be low and its measured solubility can be substantially improved  
252 by adding tags. The results were shown in Figure 4b and protein tal without tag was  
253 used as the control. All the three previously known polyionic tags increased solubility  
254 of tal when added to tal, based on experimental measurement. But none of them out-  
255 performed the tags identified in our optimization, supporting the usefulness of the tags  
256 and the optimization procedure we reported here. In addition, there was a desirable  
257 correlation between the predicted protein solubility and measured protein solubility.  
258 The linear correlation between predicted solubility and measured solubility was quan-  
259 tified by  $R^2$  with a value of 0.57. Although the previous study explored tags including  
260 aspartic acid and glutamic acid by trial and error, our study provided better optimization  
261 performance and a generally effective quantitative model.

262

## 263 **Discussions**

## 264 **Using machine learning for optimizing protein properties**

265 Using machine learning to assist the selection of proteins with specific properties has  
266 been explored recently<sup>33, 34, 35</sup>. Heckmann et al. utilized machine learning to predict the  
267 turnover number of enzymes in *E. coli* to optimize the growth rate, proteome composi-  
268 tion and physiology of organisms<sup>34</sup>. And the prediction results were further used to pa-  
269 rameterize two mechanistic genome-scale models more accurately. The machine learn-  
270 ing model was trained by using the information of protein structure, biochemistry prop-  
271 erties and assay conditions<sup>34</sup>, whereas protein sequences were used to train our predic-  
272 tion model. Therefore, their model cannot be used to optimize protein sequence for  
273 improving protein activity. Wu et al. incorporated machine learning into the directed  
274 evolution workflow to help them identify proteins with high fitness value<sup>33</sup>. Then it was  
275 applied to engineer an enzyme for stereodivergent carbon–silicon bond formation, a  
276 new-to-nature chemical transformation. However, their training data for machine learn-  
277 ing only included variants mutated at four amino acid residues. A protein might include  
278 multiple positions for mutagenesis and information of four positions is not representa-  
279 tive enough to train a machine learning model to handle other positions. The selection  
280 of mutagenesis positions need to be customized by prior knowledge on the structure of  
281 proteins. Yang et al. then reviewed the machine-learning-guided directed evolution fur-  
282 ther<sup>35</sup>. The different representation methods of protein sequence, prediction models,  
283 optimization methods, and the training data of machine learning models were discussed  
284 for different applications. Compared with the study mentioned above<sup>33, 35</sup>, we do not  
285 need to train our optimization and prediction model again when we handle a new pro-  
286 tein. In our study, we utilized the machine learning model to identify proteins with an-  
287 other desired property, protein solubility. Our training dataset was obtained by using

288 various proteins of *E. coli* and the optimization methodology did not need any custom-  
289 ization and knowledge in biochemistry for new target proteins. With only the sequence  
290 information, our optimization model can provide effective guide for improving protein  
291 solubility and activity. In addition, rather than using mutation to improve the protein  
292 properties, we added small peptide tags to improve protein solubility and activity to  
293 avoid destroying the function of the original proteins.

294

### 295 **The contribution of aspartic acid and glutamic acid**

296 In this study, we designed a novel methodology to apply a predictive model of protein  
297 solubility to improve protein solubility by adding short peptide tags. Aspartic acid and  
298 glutamic acid dominated the tags that were obtained by using our optimization strategy.  
299 This finding was consistent with the conclusion of an experiment we did to determine  
300 which amino acids were the most important in determining accuracy of our solubility-  
301 predicting model. In the experiment, we removed the percentage information of two  
302 amino acids and evaluated the negative impact on the performance of the predictive  
303 model. The model's inputs were composition of 20 amino acids, among which the per-  
304 centages of 19 amino acid were independent. As a result, removing information of only  
305 one amino acid would have no impact on model performance and we had to remove the  
306 percentages of two amino acids. We evaluated all the combinations of two amino acids.  
307 After removing aspartic acid or glutamic acid, the decrease of the prediction perfor-  
308 mance represented by  $R^2$  was the most substantial (Figure 5), indicating they were the  
309 most important ones for the model to be accurate. The causal relationship of the obser-  
310 vations from this experiment and the optimization experiment could be that these two  
311 negatively charged amino acids had large positive influence on protein solubility (as

312 seen in the optimization experiment), so they were important to the accuracy of the  
313 model prediction (as observed in the importance analysis experiment). In addition, ar-  
314 ginine which also showed some influence on the prediction performance when it was  
315 removed, did not appear in the optimization results. This might be caused by that argi-  
316 nine negatively affected the protein solubility and this hypothesis was tested (Supple-  
317 mentary Figure S4). After adding 20 arginines to the six proteins from our laboratory,  
318 all the predicted solubility was decreased. The suspected effects of glutamic acid, as-  
319 partic acid and arginine were also supported by their spearman correlation coefficients  
320 (Figure 5c), which were obtained by analyzing the large dataset we used to train our  
321 model. There were some amino acids that were identified to be important by spearman  
322 coefficient (Figure 5c) but were not found to be critical to model performance (Figure  
323 5a), such as tryptophan and phenylalanine. It may be due to that spearman coefficient  
324 alone is not sufficient to quantitatively describe the effects of amino acid on protein  
325 solubility because of its qualitative nature and it did not consider abundance of other  
326 amino acids (Figure 5b). In this study, we have shown that our machine learning model  
327 is able to quantitatively describe the relationship and guide optimization of protein se-  
328 quence.

329 When we trained the solubility-predicting model through machine learning, we did not  
330 use any biochemistry knowledge. The optimization of protein tag to maximize protein  
331 solubility was also purely mathematical without any dependence on prior knowledge.  
332 Yet, the identified most beneficial amino acids and their influence on protein solubility  
333 can be explained by using known biochemistry knowledge (electrostatic repulsion). As  
334 to why the best tags were dominated by negatively charged amino acids rather than  
335 positively charged ones, the reason might be that positively charged amino acids may  
336 also improve protein solubility but their influence is less than those of negatively

337 charged amino acids. When the number of the negatively charged amino acids was con-  
338 strained, the optimization algorithm used positively charged amino acid (lysine) to im-  
339 prove protein solubility, which led to less improvement in solubility than using the neg-  
340 atively charged ones (Supplementary Table 8 and 9).

341

## 342 **Methods**

343 **Protein database.** All the information of protein solubility used in our study is from  
344 the eSol database<sup>11</sup> which is a unique database containing continuous values of protein  
345 solubility. After removing items without sequence information according to the previ-  
346 ous study<sup>26</sup>, 3,148 proteins in the eSol dataset were used for this study. In the study  
347 which generated the dataset, the values of protein solubility were measured by synthe-  
348 sizing the recombinant proteins by cell-free protein expression technology and then  
349 being separated into soluble and insoluble fractions with centrifugation<sup>11</sup>. Solubility  
350 was defined as the ratio of supernatant protein to total protein which was quantified by  
351 SDS-PAGE.

352

353 **Training flowsheet.** The whole process of rationally engineering proteins with higher  
354 solubility includes data pre-processing, training the SVM prediction model, construct-  
355 ing an optimization methodology, and validating the methodology. As the first step,  
356 amino acid composition was extracted from protein sequences by using Amino Acid  
357 Composition Descriptor in protr package<sup>36</sup> within R software, which converted charac-  
358 ters of amino acids into numerical values indicating amino acid composition. For the  
359 second part, the SVM model was built in MATLAB and trained following the same  
360 procedure described in the previous study<sup>26</sup>. Then SVM was trained with the whole

361 dataset to predict continuous values of protein solubility from amino acid composition.  
362 For the third step, we filtered out a total of 58 proteins with low solubility of value 0.1  
363 in the original dataset and 58 proteins were picked out. Proteins with long sequences  
364 are more challenging to synthesize in experiments, therefore the protein sequences were  
365 further filtered to have less than 333.3 amino acids (1kb), which excluded 27 proteins  
366 from the eSol database. Among the 27 proteins, the one with the minimum difference  
367 between the predicted value and the real value of protein solubility, named glcE, was  
368 selected as the sample protein to build a methodology for further optimizing protein  
369 solubility. Genetic algorithm (GA), an optimization method, was explored to search for  
370 maximum predicted solubility with constraints for the sample protein. The difference  
371 between protein solubility before and after mutagenesis was used to evaluate the opti-  
372 mization effect on protein solubility. Moreover, besides the sample protein, 10 proteins  
373 with solubility of value 0.1 which have the least differences between predicted and  
374 original solubility among the 27 proteins mentioned above were selected for applying  
375 the optimization methodology. Six proteins commonly used in our laboratory were also  
376 investigated for the optimization of protein solubility. Finally, among the 16 proteins  
377 selected for optimization, 4 proteins that bear low solubility before adding the tags and  
378 high predicted solubility after adding the tags were chosen for experimental validation.  
379 The original and mutated protein sequences were synthesized to validate the change of  
380 protein solubility by measuring the protein solubility with SDS-PAGE.

381

382 **Machine learning models.** The regression version of SVM used in this study could  
383 also be named support vector regression (SVR)<sup>37</sup>. The aim of SVR is to solve<sup>38</sup>

384 minimize  $\frac{1}{2} ||w||^2$



385           subject to  $\begin{cases} y_i - \langle w, x_i \rangle - b \leq \varepsilon, \\ \langle w, x_i \rangle + b - y_i \leq \varepsilon, \end{cases}$

386   where  $x_i$  is a training sample with target value  $y_i$  and  $w$  is the normal vector to the hyperplane. The inner product plus intercept  $\langle w, x_i \rangle + b$  is the prediction value for that  
387   perplane. The difference of predicted values and true values for targets have to be within  
388   sample. The difference of predicted values and true values for targets have to be within  
389   an  $\varepsilon$  range, which is a parameter serving as a threshold.

390   A regression machine learning model SVM in MATLAB was used for optimizing protein solubility for the all the proteins in our study and was validated by experiments  
391   (Supplementary Table S10). The improved SVM model was used to optimize all the  
392   proteins *in silico* and compared with the previous one in the Discussion.  
393

394

395   **Optimization algorithms.** Genetic algorithm (GA), one of the evolutionary algorithms, is inspired by the process of natural selection observed in nature<sup>39</sup>. It is a frequently utilized randomized optimization algorithm for searching optima with constrained conditions. GA essentially simulates the way in which life evolves to find solutions to real world problems. In GA, the solutions to a problem are represented as a  
400   population of chromosomes evolving through successive generations. The offspring  
401   chromosomes are generated by merging two parent chromosomes by crossover or modifying a chromosome by mutation. The offspring chromosomes are evaluated according  
402   to the fitness or objective function in each generation. Chromosomes with higher fitness  
403   values have higher possibility to survive and the process will stop when the offspring  
404   chromosomes are almost identical or the terminal conditions set are reached. Strong  
405   individuals will dominate the generation through many iterations in the process with  
406   mutation, crossover and selection. The final chromosome represents an optimal or near-optimal solution for the optimization problem. In our problem, the chromosomes are  
407  
408

409 the sequence of peptide tags and the fitness function is the predicted solubility for pro-  
410 teins after adding tags. Several hyperparameters can be tuned for the optimization al-  
411 gorithm, such as the population size, the number of iterations for evolution and the  
412 number of individuals mutating in each generation. We used a MATLAB Toolbox to  
413 implement the optimization (iteration number = 1,000, other parameters are provided  
414 in Supplementary Table S1). The generic structure of GA in our study can be described  
415 as follows:

416 **begin:**

417       initiate a tag representing by a 20-dimensional vector with constrained condi-  
418       tions (sum of the vector is 20 and the value of each dimension is within range  
419       0-20);

420       evaluate the protein sequence after adding the tag;

421       **while** (if termination conditions are not met):

422           do crossover and mutate parent tag sequences to yield offspring se-  
423           quences;

424           evaluate the protein solubility for the proteins with offspring sequences;

425           select and generate offspring sequence with higher solubility;

426       **end while;**

427       **end.**

428

429 **Data visualization:** The heat map was plotted by using the `cmap` function of the `mat-`  
430 `plotlib` package in Python with the values of  $R^2$  after removing the information of two

431 types of amino acids. The violin plot of the amino acid compositions was made by using  
432 the violinplot function of the seaborn package in Python. Violin plot featured a kernel  
433 density estimation of the underlying distribution. Spearman's rank correlation between  
434 amino acid composition and solubility was computed using the spearmanr function of  
435 the scipy.stats package in Python. The equation used was

$$436 \rho_{spearman} = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2 \sum_i (y_i - \bar{y})^2}},$$

437 where the subscript  $i$  denoted the ranks, and  $x$  and  $y$  represented amino acid composition  
438 and solubility respectively.

439

440 **Chemicals in experimental validation:** All chemicals were purchased from Sigma-  
441 Aldrich unless otherwise stated. All reagents used were of analytical grade. The DNA  
442 oligomers used in this study were synthesized from Integrated DNA Technologies.

443

444 **Plasmid construction:** All the plasmids used in this work were constructed by using  
445 GT DNA standard<sup>40</sup> (Supplementary Table S7).

446

447 **Cell culture and SDS-PAGE analysis of protein solubility:** Each of constructed plas-  
448 mid was introduced into *E. coli* BL21 (DE3) (C2530H, New England Biolabs) for SDS-  
449 PAGE analysis by using standard heat shock protocol. In order to test the resulting  
450 strains, single colony was inoculated into 1 mL of LB with 100 µg/mL of ampicillin,  
451 and was cultured overnight at 37 °C/250 rpm. Fifty microliters of the overnight grown  
452 cell suspension were inoculated into 5 mL of K3 medium<sup>40</sup> with 100 µg/mL of ampi-  
453 cillin. When cell was grown to 0.4-0.6 optical density (OD) at 600, isopropyl β-D-1-

454 thiogalactopyranoside (IPTG) was added to a final concentration of 0.1 mM. After in-  
455 cubated overnight at 30 °C/250 rpm, the cell culture broth was diluted to OD<sub>600</sub> = 2.0,  
456 and centrifuged at 5000 g, 10 min. The obtained cell pellets were resuspended in 100  
457 µL B-PER II reagent (78248, Thermo Fisher Scientific). The mixtures were incubated  
458 for 15 min at room temperature with gentle rocking, and centrifuged at 16000 g for 20  
459 min. The obtained supernatant contained soluble cell lysates. The insoluble cell pellets  
460 were resuspended in 100 µL of 2 % w/v SDS. Both soluble and insoluble cell pellets  
461 were analyzed by using SDS-PAGE (Mini-PROTEAN® TGX™ Precast Protein Gels,  
462 4561083, Bio-Rad). The image of the gel was processed and quantified by Gel Doc EZ  
463 Gel Documentation System (Bio-Rad).

464

465 **Tal activity assay *in vitro*:** One milliliter of obtained supernatant containing soluble  
466 cell lysates was added to 4 mL of PBS buffer (pH=9.0) with 1 g/L tyrosine (final con-  
467 centration) in 50 mL falcon tube and incubated at 30 °C/250 rpm. Three hundred mi-  
468 croliters of samples were taken at 0 h, 1 h, 3 h and 12 h after incubation, and mixed  
469 with 700 µL of acetonitrile to dissolve the produced *p*-courmaric acid (PCA). The mix-  
470 ture was incubated at 30 °C/250 rpm for 1 h, and then centrifuged at 13,500 g for 5 min.  
471 Two microliters of the obtained supernatant was analyzed by using HPLC (Agilent  
472 1260 Infinity HPLC) based on a previously reported method<sup>40</sup>.

473

#### 474 **Supplementary information**

475 Supplementary data are available online.

476

## 477 **Codes availability**

478 We present the optimization workflow as a series of notebooks hosted on GitHub  
479 ([https://github.com/xiaomizhou616/optimization\\_protein-solubility](https://github.com/xiaomizhou616/optimization_protein-solubility)). The workflow  
480 can be used as a template for analysis of other expression and solubility datasets.

481

## 482 **Reference**

- 483 1. Itakura K, *et al.* Expression in *Escherichia coli* of a chemically synthesized gene for the hormone  
484 somatostatin. *Science* **198**, 1056-1063 (1977).  
485
- 486 2. Chan W-C, Liang P-H, Shih Y-P, Yang U-C, Lin W-c, Hsu C-N. Learning to predict expression  
487 efficacy of vectors in recombinant protein production. *BMC Bioinform* **11**, S21 (2010).  
488
- 489 3. Fang H, Li D, Kang J, Jiang P, Sun J, Zhang D. Metabolic engineering of *Escherichia coli* for de  
490 novo biosynthesis of vitamin B 12. *Nature communications* **9**, 4917 (2018).  
491
- 492 4. Lempp M, Farke N, Kuntz M, Freibert SA, Lill R, Link H. Systematic identification of metabolites  
493 controlling gene expression in *E. coli*. *Nature communications* **10**, 1-9 (2019).  
494
- 495 5. Esposito D, Chatterjee DK. Enhancement of soluble protein expression through the use of  
496 fusion tags. *Current opinion in biotechnology* **17**, 353-358 (2006).  
497
- 498 6. Idicula - Thomas S, Balaji PV. Understanding the relationship between the primary structure  
499 of proteins and its propensity to be soluble on overexpression in *Escherichia coli*. *Protein Sci*  
500 **14**, 582–592 (2005).  
501
- 502 7. Magnan CN, Randall A, Baldi P. SOLpro: accurate sequence-based prediction of protein  
503 solubility. *Bioinformatics* **25**, 2200–2207 (2009).  
504
- 505 8. Trésaugues L, *et al.* Refolding strategies from inclusion bodies in a structural genomics project.  
506 *Journal of structural and functional genomics* **5**, 195-204 (2004).  
507
- 508 9. Ganesan A, *et al.* Structural hot spots for the solubility of globular proteins. *Nature*  
509 *communications* **7**, 10816 (2016).  
510
- 511 10. Reyes LH, Cardona C, Pimentel L, Rodríguez-López A, Alméjiga-Díaz CJ. Improvement in the  
512 production of the human recombinant enzyme N-acetylgalactosamine-6-sulfatase (rhGALNS)  
513 in *Escherichia coli* using synthetic biology approaches. *Scientific reports* **7**, 5844 (2017).  
514
- 515 11. Niwa T, *et al.* Bimodal protein solubility distribution revealed by an aggregation analysis of the  
516 entire ensemble of *Escherichia coli* proteins. *Proc Natl Acad Sci USA* **106**, 4201–4206 (2009).  
517
- 518 12. Zhou K, Zou R, Stephanopoulos G, Too H-P. Enhancing solubility of deoxyxylulose phosphate  
519 pathway enzymes for microbial isoprenoid production. *Microbial cell factories* **11**, 148 (2012).  
520
- 521 13. Kronqvist N, *et al.* Efficient protein production inspired by how spiders make silk. *Nature*  
522 *communications* **8**, 15504 (2017).  
523

- 524 14. Wu Y, *et al.* Google's neural machine translation system: Bridging the gap between human and  
525 machine translation. *arXiv preprint arXiv:160908144*, (2016).  
526
- 527 15. Bojarski M, *et al.* End to end learning for self-driving cars. *arXiv preprint arXiv:160407316*,  
528 (2016).  
529
- 530 16. LeCun Y, Bengio Y, Hinton G. Deep learning. *nature* **521**, 436 (2015).  
531
- 532 17. Silver D, *et al.* Mastering the game of Go with deep neural networks and tree search. *nature*  
533 **529**, 484 (2016).  
534
- 535 18. Ferrucci D, Levas A, Bagchi S, Gondek D, Mueller E. Watson: Beyond Jeopardy! *Artif Intell.*  
536 (2013).  
537
- 538 19. Godec P, *et al.* Democratized image analytics by visual programming through integration of  
539 deep models and small-scale machine learning. *Nature Communications* **10**, 1-7 (2019).  
540
- 541 20. Weber T, Wiseman NA, Kock A. Global ocean methane emissions dominated by shallow coastal  
542 waters. *Nature Communications* **10**, 1-10 (2019).  
543
- 544 21. Li L, *et al.* Machine-learning reprogrammable metasurface imager. *Nature communications* **10**,  
545 1082 (2019).  
546
- 547 22. Zador AM. A critique of pure learning and what artificial neural networks can learn from animal  
548 brains. *Nature communications* **10**, 1-7 (2019).  
549
- 550 23. Diaz AA, Tomba E, Lennarson R, Richard R, Bagajewicz MJ, Harrison RG. Prediction of protein  
551 solubility in *Escherichia coli* using logistic regression. *Biotechnol Bioeng* **105**, 374–383 (2010).  
552
- 553 24. Agostini F, Vendruscolo M, Tartaglia GG. Sequence-based prediction of protein solubility. *J Mol*  
554 *Biol* **421**, 237–241 (2012).  
555
- 556 25. Xiaohui N, Feng S, Xuehai H, Jingbo X, Nana L. Predicting the protein solubility by integrating  
557 chaos games representation and entropy in information theory. *Expert Syst Appl* **41**, 1672–  
558 1679 (2014).  
559
- 560 26. Han X, Wang X, Zhou K. Develop machine learning based regression predictive models for  
561 engineering protein solubility. *Bioinformatics*, (2019).  
562
- 563 27. Paraskevopoulou V, Falcone F. Polyionic tags as enhancers of protein solubility in recombinant  
564 protein expression. *Microorganisms* **6**, 47 (2018).  
565
- 566 28. Jendresen CB, *et al.* Highly active and specific tyrosine ammonia-lyases from diverse origins  
567 enable enhanced production of aromatic compounds in bacteria and *Saccharomyces*  
568 *cerevisiae*. *Appl Environ Microbiol* **81**, 4458-4476 (2015).  
569
- 570 29. Rodriguez A, Kildegaard KR, Li M, Borodina I, Nielsen J. Establishment of a yeast platform strain  
571 for production of p-coumaric acid through metabolic engineering of aromatic amino acid  
572 biosynthesis. *Metabolic engineering* **31**, 181-188 (2015).  
573
- 574 30. Bianchi E, Venturini S, Pessi A, Tramontano A, Sollazzo M. High level expression and rational  
575 mutagenesis of a designed protein, the minibody: from an insoluble to a soluble molecule. *J*  
576 *Mol Biol* **236**, 649-659 (1994).  
577
- 578 31. Chan P, Curtis RA, Warwicker J. Soluble expression of proteins correlates with a lack of  
579 positively-charged surface. *Scientific reports* **3**, 3333 (2013).  
580

- 581 32. Nguyen TKM, Ki MR, Son RG, Pack SP. The NT11, a novel fusion tag for enhancing protein  
582 expression in *Escherichia coli*. *Applied microbiology and biotechnology*, 1-12 (2019).  
583
- 584 33. Wu Z, Kan SJ, Lewis RD, Wittmann BJ, Arnold FH. Machine learning-assisted directed protein  
585 evolution with combinatorial libraries. *Proc Natl Acad Sci USA* **116**, 8852-8858 (2019).  
586
- 587 34. Heckmann D, *et al.* Machine learning applied to enzyme turnover numbers reveals protein  
588 structural correlates and improves metabolic models. *Nature communications* **9**, 5252 (2018).  
589
- 590 35. Yang KK, Wu Z, Arnold FH. Machine-learning-guided directed evolution for protein  
591 engineering. *Nature methods*, 1 (2019).  
592
- 593 36. Xiao N, Xu Q, Cao D. Protr: Protein sequence descriptor calculation and similarity computation  
594 with R. R package version 0.2-1. (2014).  
595
- 596 37. Drucker H, Burges CJ, Kaufman L, Smola AJ, Vapnik V. Support vector regression machines. In:  
597 *Advances in neural information processing systems* (ed<sup>^</sup>(eds) (1997).  
598
- 599 38. Smola AJ, Schölkopf B. A tutorial on support vector regression. *Statistics and computing* **14**,  
600 199-222 (2004).  
601
- 602 39. Mitchell M. An introduction to genetic algorithms mit press. *Cambridge, Massachusetts*  
603 *London, England*, (1996).  
604
- 605 40. Ma X, *et al.* A standard for near-scarless plasmid construction using reusable DNA parts. *Nature*  
606 *communications* **10**, 3294 (2019).  
607  
608

## 609 **Acknowledgments**

610 We acknowledge the MOE Research Scholarship, MOE Tier-1 grant (R-279-000-452-  
611 133) and NRF CRP grant (R-279-000-512-281) in Singapore. We thank Cortes-Pena  
612 Yoel Rene for providing data visualization for data distribution and Spearman's rank  
613 correlation tornado plot.

614

## 615 **Author contributions**

616 X.H. developed the optimization algorithms and statistical analyses. W.N. performed  
617 the experimental preparation and validation, and X.M designed and guided the experi-  
618 ments. All of them were supervised by X.W. and K.Z.. X.H. and W.N. wrote the man-  
619 uscript with inputs from all the co-authors. All authors discussed the results and com-  
620 mented on the manuscript.

621

622 *Conflict of Interest:* none declared.

623

624

625

626

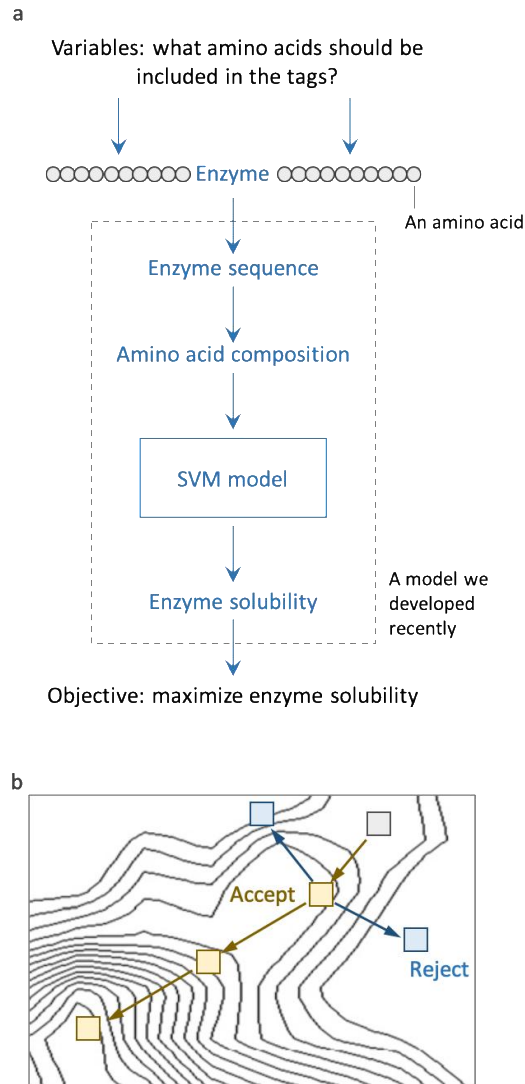
627

628

629

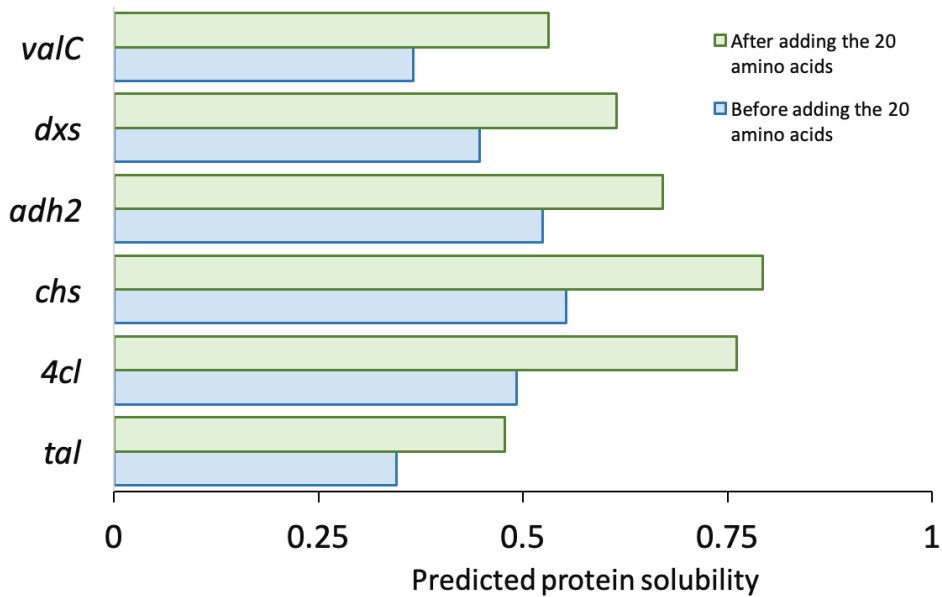
630





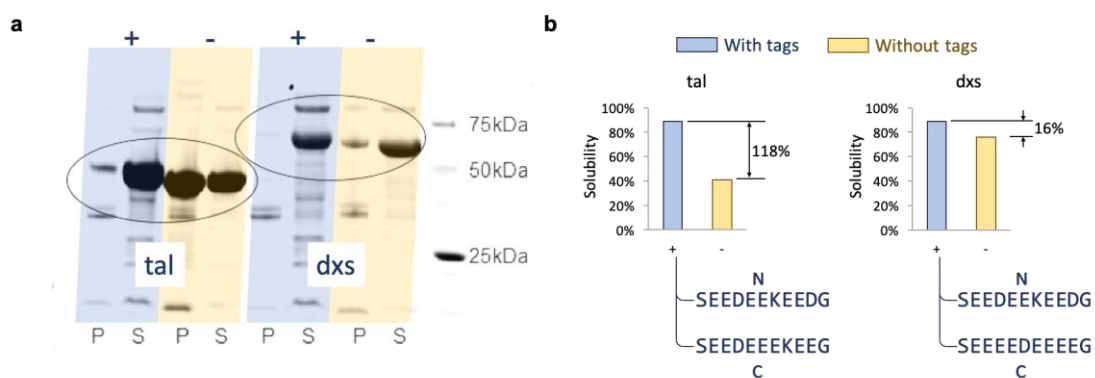
631

632 **Figure 1** Machine learning model assisted optimization of protein solubility. (a) Illustration of the ob-  
633 jective function when we aimed to improve protein solubility by adding short peptide tags. SVM: support  
634 vector machine. A SVM regression model we recently developed was used in this study<sup>26</sup>. (b) Illustration  
635 of the optimization algorithm. Genetic algorithm was used in this study.



636

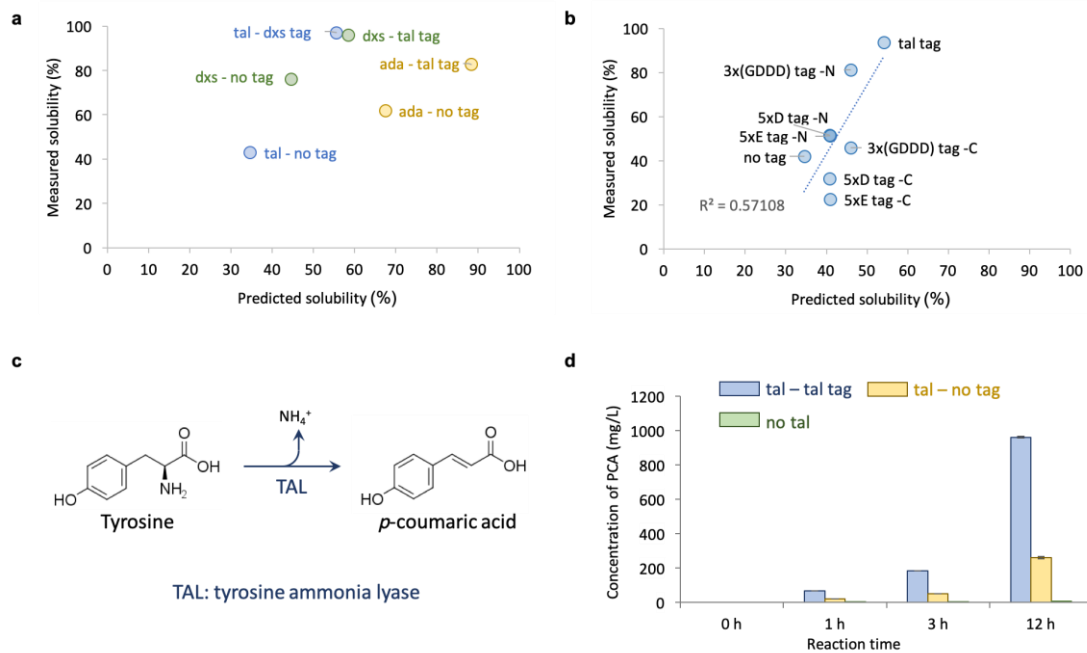
637 **Figure 2** The predicted solubility before and after adding 20 amino acids for six proteins commonly used  
 638 by our laboratory. The six proteins were valC (valencene synthase), dxs (1-deoxy-D-xylulose-5-phos-  
 639 phate synthase), adh2 (alcohol dehydrogenase), chs (chalcone synthase), 4cl (4-coumarate-CoA ligase)  
 640 and tal (tyrosine ammonia-lyase). Their sequences were listed in Supplementary Table S7. Before adding  
 641 the tags, the protein solubility of them was predicted by SVM and recorded. Then GA was used to opti-  
 642 mize their solubility by adding 20 amino acids. The protein solubility after adding the tags was also  
 643 recorded for comparison.



644

645 **Figure 3 (a)** The SDS-PAGE analysis of protein tal and dxs expressed in *E. coli* with and without tags  
 646 designed by our optimization algorithm. “+” and “-” represented expressed proteins with and without  
 647 peptide tags respectively. “P” and “S” represented the pellet fraction (insoluble) and supernatant fraction  
 648 (soluble), respectively. The oval shapes highlight the bands of dxs and tal proteins. Protein tal and dxs

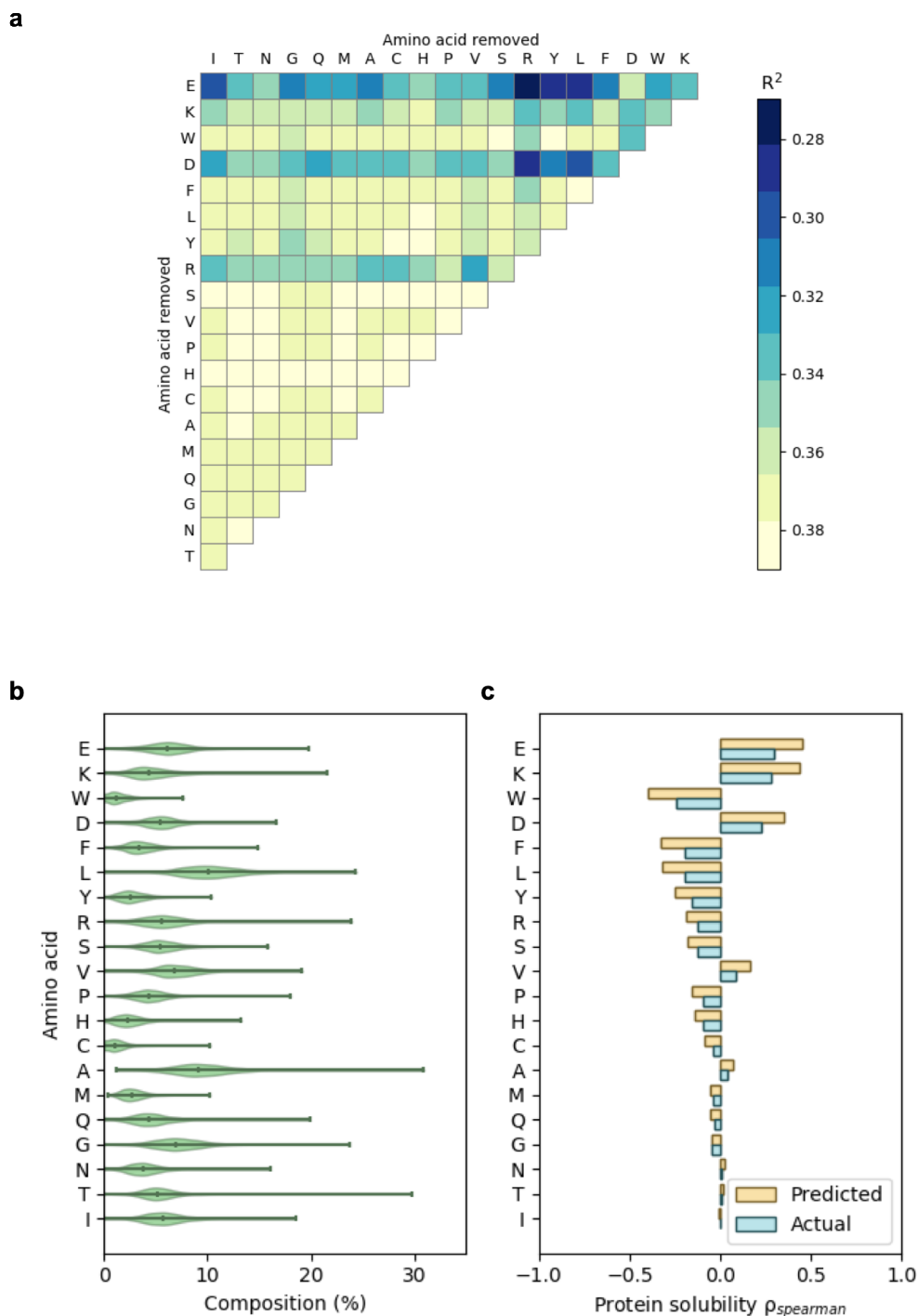
649 were expressed in K3 medium with 20 g/L glucose at 30 °C. **(b)** Quantitative presentation of the SDS-  
 650 PAGE images in **a**. The protein solubility was the ratio of soluble protein amount to the total protein  
 651 amount. The protein amount was estimated by using band intensity on SDS-PAGE images. The se-  
 652 quences of the designed tags for N-terminal and C-terminal were shown. The amino acid S and G on the  
 653 two ends of the tags were the linkers for GT DNA assembly standard, which was used to construct the  
 654 plasmids in this study<sup>40</sup>.



655  
 656 **Figure 4 (a)** The predicted and measured solubility of tal, dxs and ada after adding tags designed for  
 657 other proteins. The purpose of switching tags for proteins was to test if the solubility-enhancing tags are  
 658 generally effective in improving protein solubility. The same protein was labelled by using the same  
 659 color to highlight the data before and after adding tags. In the data labels, the text before “-” indicates  
 660 protein name and the text after “-” indicates the tags used if any. In the process of measuring the solubil-  
 661 ity, the protein expression condition was K3 medium with 20 g/L glucose at 30 °C. **(b)** The comparison  
 662 of the tags designed in this study with tags used in previous studies. Protein tal was the only model  
 663 protein used in this plot. No tag: solubility of tal without any tag. Tal tag: solubility of tal when we added  
 664 the tags that were designed by our optimization algorithm for tal. 5xE tag -N/C: solubility of tal when  
 665 5xE tag (EEEEEE) was added to its N- or C-terminus. 5xD tag -N/C: solubility of tal when 5xD tag  
 666 (DDDDD) was added to its N- or C-terminus. 3x(GDDD) -N/C: solubility of tal when 3x(GDDD) tag  
 667 (GDDDGDDDGDDD) was added to its N- or C-terminus. 5xD, 5xE and 3x(GDDD) were three tags  
 668 used in a previous study and used here for comparison<sup>27</sup>. Since in previous study, only one tag was added

669 to one protein, either at N- or C-terminus, we tested both cases for each tag. The two tags we designed  
670 for tal were added to both ends of tal (Figure 1 and 3b). The sequences of all the tags are provided in  
671 Supplementary Table S7. In the process of measuring the solubility, the protein expression condition was  
672 K3 medium with 20 g/L glucose at 30 °C. (c) The reaction catalyzed by enzyme tal. (d) The protein  
673 activity of protein tal before and after introducing tal tag. The product of the reaction catalyzed by en-  
674 zyme tal was p-coumaric acid (PCA) and its concentration was used to indicate the activity of protein  
675 tal. Cell lysate containing tal was used in the reaction. tal – tal tag: the strain containing tal with the tags  
676 designed in this study. Tal – no tag: the strain containing tal without any tag. No tal: the strain that did  
677 not express tal. The bars indicate the mean of six replicates. The error bars indicate standard error of six  
678 replicates.

679



680

681 **Figure 5 (a)** Importance of various amino acids in determining the accuracy of the SVM regression  
 682 model. The  $R^2$  of the SVM model was shown by using a heat map after removing the information of two  
 683 types of amino acids. Model training is described in Materials and Methods. Single letter amino acid  
 684 abbreviations are used in this figure. All the combinations of removing two types of amino acids are  
 685 tested and the performance of the resulting models is presented in the upper triangular matrix. Perfor-  
 686 mance of the models was gauged by using  $R^2$ , which is presented here by using color (a color bar is

687 provided). The darker the color is, the more important the related amino acids are to the model perfor-  
688 mance. **(b)** The distribution of amino acid composition (the input variables of the SVM model we used)  
689 among all the proteins in the eSol database (the data source we used to train the SVM model). The violin  
690 plot showed the mean value and the range of the amino acid composition used to train the SVM model.  
691 **(c)** The Spearman's rank correlation between actual/predicted protein solubility and various amino acids.  
692 Spearman's correlation,  $\rho_{spearman}$ , is a measure of monotonicity and represents the general sensitivity  
693 of solubility to amino acid composition. A comparison between the Spearman's rank correlation tornado  
694 plot for actual solubility and predicted solubility depicted how the model captured and magnified general  
695 trends between amino acid composition and solubility. For example, for both the actual and predicted  
696 solubility of proteins in the eSol dataset, the composition of D, E, or K was positively correlated with  
697 solubility.  
698