

# 1 **Benchmark and optimization of AlphaFold structures based** 2 **virtual screening strategy**

3 Yanfei Peng<sup>1</sup>, Xia Wu<sup>1</sup>, Liang Lin<sup>1</sup>, Zhiluo Deng<sup>2\*</sup>, Limin Zhao<sup>1\*</sup>, Hao Ke<sup>1\*</sup>

4

5 <sup>1</sup> Human Aging Research Institute (HARI) and School of Life Science, Nanchang  
6 University, and Jiangxi Key Laboratory of Human Aging, Nanchang, 330031 Jiangxi  
7 China

8 <sup>2</sup> Department for Computational Biology of Infection Research, Helmholtz Center for  
9 Infection Research, Braunschweig, Germany; Braunschweig Integrated Centre of  
10 Systems Biology (BRICS), Technische Universität Braunschweig, Braunschweig,  
11 Germany;

12

13 \* Corresponding authors:

14 Zhiluo Deng, Tel: +49 531-391-55273, E-mail: Zhiluo.Deng@helmholtz-hzi.de

15 Limin Zhao, Tel: +86 181-8600-5898, E-mail: zhaolimin@ncu.edu.cn

16 Hao Ke, Tel: +86 135-1874-2153, E-mail: kehao@ncu.edu.cn

17

18

## 19 **Abstract**

20 Recent advancements in artificial intelligence such as AlphaFold, have enabled more  
21 accurate prediction of protein three-dimensional structure from amino acid sequences.  
22 This has attracted significant attention, especially for the application of AlphaFold in  
23 drug discovery. Moreover, how to take full advantage of AlphaFold to assist with  
24 virtual screening remains elusive. We comprehensively evaluated the AlphaFold  
25 structures of 51 selected targets from the DUD-E database in virtual screening. Our

26 analyses show that the virtual screening performance of about 35% of the AlphaFold  
27 structures was equivalent to that of DUD-E structures, and about 25% of the  
28 AlphaFold structures yielded better results than the DUD-E structures. Remarkably,  
29 for the 23 targets, AlphaFold structures produced slightly better results than the Apo  
30 structures. Moreover, we developed a new consensus scoring method based on Z-  
31 score standardization and exponential function, which showed improved screening  
32 performance compared to traditional scoring methods. By implementing a multi-stage  
33 virtual screening process and the new consensus scoring method, we were able to  
34 improve the speed of virtual screening by about nine times without compromising the  
35 enrichment factor. Overall, our results provide insights into the potential use of  
36 AlphaFold in drug discovery and highlight the value of consensus scoring and multi-  
37 stage virtual screening.

38

39 **Key words:** AlphaFold, virtual screening, high-throughput docking, consensus  
40 scoring, multi-stage virtual screening

41

## 42 **Introduction**

43 AlphaFold is deep learning based program that can predict the three-dimensional  
44 structure of proteins from amino acid sequences. It is one of the low-cost computing  
45 methods to obtain highly accurate protein structures<sup>1</sup>. Currently, the AlphaFold  
46 protein structure database contains over 200 million entries, covering the human  
47 proteome and the proteome of 47 other important organisms crucial to research and  
48 global health. The availability of large numbers of easily accessible and highly  
49 accurate protein structures has made AlphaFold a valuable resource for research fields  
50 related to protein structure, particularly for protein research that lacks experimental  
51 structure information<sup>2</sup>.

52 Several studies have explored the use of AlphaFold protein structure for drug  
53 discovery<sup>3-5</sup>. A few recent studies evaluated the performance of AlphaFold protein

54 structure in virtual screening. Wong et al. evaluated the effectiveness of AlphaFold  
55 protein structure to predict the binding affinity of 296 E. coli proteins and 218 small  
56 molecules with known antibacterial activity, and compared the results to the  
57 experimental structure of 12 proteins<sup>6</sup>. Scardino et al. assessed the performance of the  
58 structure in the PDB database and AlphaFold structure in virtual screening using four  
59 docking software and ECR (Exponential Consensus Ranking, based on ranking) on 16  
60 protein targets, and the consensus score of PRC (Pose/Ranking Consensus, based on  
61 docking pose and ranking)<sup>7</sup>. Zhang examined the efficiency of Holo, Apo and  
62 AlphaFold structures in virtual screening for 28 targets et al. with Glide molecular  
63 docking method<sup>8</sup>. Alon, et al. found that although the AlphaFold structure of  $\sigma$ 2  
64 receptor is very similar to the crystal structure, the score of small molecule  
65 compounds with the AlphaFold structure is lower than that of the crystal structure<sup>9</sup>.

66 Virtual screening based on molecular docking is a widely used method in  
67 computer-aided drug design. With scoring function, molecular docking can effectively  
68 identify small molecules that interact with the ligand binding pocket of receptor  
69 protein. Numerous studies have reported the use of molecular docking based virtual  
70 screening to discover inhibitors and agonists of target proteins<sup>9-12</sup>. Virtual screening  
71 can effectively reduce the cost of drug discovery and accelerate the process. There are  
72 many different molecular docking methods widely used in virtual screening, including  
73 AutoDock Vina<sup>13,14</sup>, Qvina2<sup>15</sup>, idock<sup>16</sup>, ICM<sup>17</sup>, Glide<sup>18</sup>, Gold<sup>19</sup>, etc. However,  
74 molecular docking still has high false positive rate which limits its use in drug  
75 discovery. Some recent studies have used the consensus scoring method to integrate  
76 the score reported by different docking software to reduce the false positive rate and  
77 improve the enrichment factor<sup>20,21</sup>. Other studies have even takes into account the  
78 docking pose in consensus scoring, which can further improve the performance of the  
79 virtual screening<sup>22-24</sup>. Therefore, further optimization of the consensus scoring  
80 method could be of great significance for improving the value of AlphaFold in virtual  
81 screening.

82 With the rapid expansion of the small molecule library for virtual screening, for  
83 example, the ZINC database has provided more than 120 million purchasable drug-

84 like compounds<sup>25</sup>, large-scale virtual screening plays an increasingly important role in  
85 drug discovery. However large-scale virtual screening requires a substantial amount of  
86 computing resources. To reduce the runtime and save computing resources, the multi-  
87 stage screening method has been used in large-scale virtual screening. Gorculla et al.  
88 conducted a large-scale screening with 1.3 billion compounds using a multi-stage  
89 screening method. In the first stage, Qvina2 was used for rapid screening with a low  
90 accuracy. In the second stage, 13 residues of the receptor were considered flexible,  
91 and AutoDock Vina and Smina Vinardo were used to re-score the top 3 million  
92 compounds from the first stage with a higher accuracy<sup>10</sup>.

93 Protein structure can be categorized into Holo structure and Apo structure based  
94 on the presence of small molecules in the binding pocket. Holo structures have small  
95 molecules in the binding pocket, while Apo structures do not. Previous studies have  
96 shown that the category of the protein structure has a significant impact on virtual  
97 screening, and the results of Holo structure is generally better than that of Apo  
98 structure<sup>26,27</sup>.

99 To comprehensively evaluate the impact of using AlphaFold protein structures,  
100 Holo and Apo structures on virtual screening, we selected 51 protein targets from the  
101 widely used DUD-E database<sup>28</sup> and more than 400,000 small molecules. A novel  
102 consensus scoring function was developed to achieve a lower false discovery rate. In  
103 the benchmark, various molecular docking software, scoring functions, and consensus  
104 scoring including the novel consensus scoring method were used. Our analyses show  
105 that the virtual screening performance of about 35% of the AlphaFold structures was  
106 equivalent to that of DUD-E structures, and about 25% of the AlphaFold structures  
107 yielded better results than the DUD-E structures. Remarkably, for the 23 targets,  
108 AlphaFold structures produced slightly better results than the Apo structures. Notably,  
109 the new consensus scoring method is superior to the traditional scoring method and  
110 can be used to design a multi-stage virtual screening process that improves the virtual  
111 screening speed by about 9 times without compromising the enrichment factor EF1%.

112

113

## 114 **Results**

### 115 **Virtual screening based on protein structure in DUD-E database and AlphaFold**

116 With autodock4 scoring function, all the metrics did not show substantial  
117 difference between DUD-E protein structures and AlphaFold protein structures (Fig.  
118 1b, Fig. 1c) (p-value for logAUC between DUD-E and AlphaFold is 0.081). However,  
119 with the other scoring functions, DUD-E structures performed better than AlphaFold  
120 structures overall in terms of all metrics except for AUC (Fig. 1c). For idock the EF1%  
121 was  $6.342 \pm 8.808$  for DUD-E structure and  $3.985 \pm 5.191$  for AlphaFold structure  
122 (p-value = 0.010). The logAUC based on idock for DUD-E structure was  $0.252 \pm$   
123  $0.113$  while  $0.228 \pm 0.088$  for AlphaFold structure (p-value = 0.021). Notably, the  
124 virtual screening based on autodock4 scoring function yielded better results on all  
125 metrics. The results suggest that the protein structure in DUD-E database was slightly  
126 better than AlphaFold predicted structure for virtual screening. The selection of  
127 scoring functions had a great impact on the results of virtual screening (Table 1). The  
128 results for autodock4, idock, rf\_score and vinardo based on DUD-E structures and  
129 AlphaFold structures of every target are shown in Supplementary Information Table  
130 2-9.

131

### 132 **Virtual screening results of Holo, AlphaFold and Apo protein structures**

133 Previous studies have shown that the virtual screening result of Holo structure is  
134 generally better than that of Apo structure. The protein structure in the DUD-E  
135 database is Holo structure. To evaluate whether AlphaFold protein structure produces  
136 better result compared to Apo structure in virtual screening, we tried to collect the  
137 Apo structure for the 51 selected targets in DUD-E database. However, only 23 out of  
138 these 51 targets have Apo structure available in the PDB database with known binding  
139 pocket. We then performed virtual screening with the Apo structures of these 23  
140 targets.

141 Molecular docking is a technique that involves fitting small molecules into the  
142 binding pocket of a protein structure and allows for the identification of appropriate  
143 conformations. Therefore, the binding pocket is essential for virtual screening using  
144 molecular docking. When analyzing the protein structures, we found that HIVRT has  
145 clear binding pocket in its Holo and Apo structures encompassing the small co-  
146 crystalline molecules. However, no apparent binding pocket was found in its  
147 AlphaFold structure (Fig. 2a), so HIVRT was not included in the 23 targets mentioned  
148 above. On the other hand, AKT2, has clear binding pockets in its Holo and AlphaFold  
149 structures but not in the Apo structure (Fig. 2b). Therefore, AKT2 was also excluded  
150 from the selected targets.

151 The AUC values between Holo and Apo groups were significantly different  
152 when using scoring function autodock4 (Holo  $0.628 \pm 0.154$ , Apo  $0.580 \pm 0.142$ , p-  
153 value = 0.040), idock (Holo  $0.641 \pm 0.152$ , Apo  $0.603 \pm 0.144$ , p-value = 0.043).  
154 With the vinardo function, AlphaFold structures were significantly better than Apo  
155 structures in terms of EF1% (AlphaFold  $4.183 \pm 4.919$ , Apo  $2.311 \pm 2.455$ , p-value  
156 = 0.032) and logAUC (AlphaFold  $0.230 \pm 0.087$ , Apo  $0.210 \pm 0.074$ , p-value =  
157 0.042). It shows that the protein structure in DUD-E database is better than Apo  
158 structure in virtual screening. The protein structure in DUD-E database is slightly  
159 better than AlphaFold protein structure for virtual screening, and AlphaFold protein  
160 structure is slightly better than Apo protein structure in virtual screening.

161 It is worth noting that the difference between DUD-E and AlphaFold groups in  
162 Figure 2c is relatively smaller than the difference in Figure 1c. It is because the  
163 number of protein targets tested in Figure 2c is less. The virtual screening of protein  
164 structure in DUD-E database (the average of logAUC for idock is 0.244) was better  
165 than the average value of AlphaFold protein structure (the average of logAUC for  
166 idock is 0.220) in different scoring functions and different indicators. And the virtual  
167 screening of AlphaFold protein structure is better than the overall scores of protein  
168 Apo structure (the average of logAUC for idock is 0.202), while different scoring  
169 functions have a significant impact on the results of virtual screening (Table 2). The

170 results for autodock4, idock, rf\_score and vinardo based on Apo structures of every  
171 target are shown in Supplementary Information Table 10-13.

172

### 173 **Evaluation of consensus scoring method based on score or ranking**

174 Previous studies have shown that consensus scoring based on multiple scoring  
175 functions<sup>20-23</sup> can take advantages of each scoring function to achieve better screening  
176 results. However, different scoring functions may need to be standardized before they  
177 can be integrated. Common standardization methods include ranking, AASS, Average  
178 of auto scaled scores, Z-score scaling.

179 In this study, we selected four appropriate scoring functions (autodock4, idock,  
180 rf\_score, and vinardo) and combined them using 12 different consensus scoring  
181 methods (3 standardization methods  $\times$  4 consensus calculation methods). We tested  
182 these methods on 51 previously selected targets, with the results of the four single  
183 scoring functions as controls.

184 Different scoring functions may produce similar results, which will not help  
185 improving the consensus outcome. So, we tested the correlation between each one  
186 another among the five scoring functions including autodock4, idock, qvina, rf\_score  
187 and vinardo. The correlation coefficient  $R^2$  between qvina and idock is greater than  
188 0.9 (for logAUC is 0.937, for BEDROC( $\alpha=80.5$ ) is 0.930) and the fitting line (for  
189 logAUC is  $y = 0.909x + 0.018$ , for BEDROC( $\alpha=80.5$ ) is  $y = 1.058x + 0.002$ ) is  
190 close to  $y=x$  (Fig. 3a), while the correlation coefficient  $R^2$  between other pairs are less  
191 than 0.75. These results indicated that qvina and idock have a high similarity. In fact,  
192 qvina and idock are both developed based on AutoDock Vina. Therefore, we finally  
193 selected autodock4, idock, rf\_score and vinardo scoring functions for consensus  
194 scoring.

195 To show four different calculation methods of consensus scores, we compared  
196 the distribution score of idock, rf\_score and consensus evaluation methods on the  
197 representative BACE1 protein (Fig. 3b). From the distribution of active compounds  
198 and decoy compounds, the top small molecules were more likely to be active

199 compounds (for idock the EF10% is 2.66 while for rf\_score the EF10% is 2.87),  
200 indicating that idock and rf\_score has certain ability to enrich active compounds.  
201 Moreover, we found that the best exp\_z\_score, consensus score of exponential  
202 function based on Z-score scaling, was generally better than that of autodock4, idock  
203 and rf\_score and vinardo scoring functions. The exp\_z\_score were also superior to  
204 most other consensus scoring methods (Fig. 3c).

205 The exp\_rank (ECR) was also a consensus scoring method based on ranking  
206 exponential function. Compared the consensus score of ECR and exp\_z\_score groups  
207 by paired t-test, exp\_z\_score showed better than ECR method (for example, logAUC  
208 for exp\_z\_score is 0.290 while for ECR is 0.284), except for the AUC (for  
209 exp\_z\_score is 0.686 while for ECR is 0.689). Notely, BEDROC( $\alpha=20.0$ ) and EF10%  
210 of exp\_z\_score (0.253 and 3.202) are higher than that of ECR (0.244 and 3.104) with  
211 significant difference (p-value is 0.028 and 0.028)(Table 3). The results for scoring  
212 function and consensus scoring function (including using 4, 3, 2 scoring functions)  
213 based on DUD-E structures, AlphaFold structures and Apo structures are shown in  
214 Supplementary Information Table 14-16.

215

### 216 **Re-evaluation of Holo, AlphaFold and Apo protein structures using consensus** 217 **scoring method**

218 We also evaluated the virtual screening results of Holo, AlphaFold and Apo  
219 protein structures with four representative consensus scores. Again, with the  
220 consensus scoring functions, DUD-E structure generated better results compared to  
221 AlphaFold protein structure (such as for exp\_z\_score, the logAUC for DUD-E is  
222  $0.290 \pm 0.133$  while the logAUC for AlphaFold is  $0.256 \pm 0.096$ , and the p-value  
223 for logAUC between DUD-E and AlphaFold is 0.014)(Fig. 4a). The average values of  
224 main data groups have substantial differences, while the direct average values of a  
225 few data groups do not have significant differences. Remarkably, AlphaFold structure  
226 is slightly better than Apo structure in virtual screening with the consensus scores  
227 (such as for exp\_z\_score, the EF1% for AlphaFold is  $5.161 \pm 5.284$  while the EF1%



228 for Apo is  $3.223 \pm 5.321$ , and the p-value for EF1% between AlphaFold and Apo is  
229 0.050)(Fig. 4b).

230 In addition, for each of the 51 previously selected targets, we also made a  
231 difference in each index of the AlphaFold protein structure and protein Holo structure  
232 screening results, namely AlphaFold – DUD-E (Fig. 4c). If the difference of defined  
233 logAUC is less than -0.03, it indicates that the virtual screening result of AlphaFold  
234 protein structure is worse than that of DUD-E protein structure. If the difference of  
235 defined logAUC is greater than or equal to -0.03 and less than or equal to 0.03, it  
236 indicates that the virtual screening result of AlphaFold protein structure is equivalent  
237 to that of DUD-E protein structure. If the difference of defined logAUC is greater  
238 than 0.03, it indicates that the virtual screening result of AlphaFold protein structure is  
239 better than that of DUD-E protein structure. There are 20 targets (about 39%) with  
240 AlphaFold protein structure worse than DUD-E protein structure, 18 targets (about  
241 35%) with AlphaFold protein structure equivalent to DUD-E protein structure, and 13  
242 targets (about 25%) with AlphaFold protein structure better than DUD-E protein  
243 structure.

244 We also studied the 51 protein Holo structure and AlphaFold protein structure for  
245 `exp_z_score` by paired t-test. The average AlphaFold minus DUD-E of all indicators  
246 were negative value, and there is a significant difference between the average  
247 AlphaFold and DUD-E of most indicators (Table 4). These similar results were also  
248 observed in the situation of that of Apo minus AlphaFold (Table 5). With consensus  
249 scoring of virtual screening, Holo protein structures from DUD-E database were  
250 better than AlphaFold protein structure, which were better than Apo structures. The  
251 results for `exp_z_score` consensus scoring function based on DUD-E structures,  
252 AlphaFold structures and Apo structures of every target are shown in Supplementary  
253 Information Table 17-19.

254

255

256 **Multi-stage screening combined with consensus scoring**

257        Although the virtual screening effect of the best consensus scoring obtained in  
258 the test is better than that of a single scoring function, consensus scoring is based on  
259 multiple scoring functions. The computing resources consumed in the whole process  
260 of consensus scoring are the sum of all the computing resources consumed by the  
261 involved single scoring function and docking software, which consumes more  
262 computing resources.

263        With the rapid expansion of the small molecule compound library that can be  
264 used for virtual screening, large-scale virtual screening is also increasingly widely  
265 used. Therefore, it is of great significance to improve the efficiency of virtual  
266 screening and reduce the computing resources consumed by virtual screening to save  
267 the cost of large-scale virtual screening and accelerate the drug screening process.  
268 Several studies have used multi-stage screening to improve the efficiency of virtual  
269 screening. Multi-stage screening uses rough but fast screening in the first part of the  
270 stage, and more refined but slow screening in the second part of the stage, so as to  
271 give consideration to the early enrichment ability and computational efficiency of  
272 virtual screening.

273        In order to use the advantages of consensus scoring to improve the ability of  
274 virtual screening and reduce the computational resources consumed in the whole  
275 process of consensus scoring, we propose a multi-stage screening combined with  
276 consensus scoring (Fig. 5a). It is worth noting that since our study above has  
277 performed virtual screening on 51 targets, and the scores of different scoring functions  
278 between all proteins and small molecules are known, multi-stage screening (plan A, B,  
279 C, D, E, F) does not actually perform virtual screening, but uses corresponding known  
280 data, and the time consumed is also a theoretical estimate calculated based on known  
281 data.

282        On the crucial EF1% index for evaluating the early enrichment ability of virtual  
283 screening, the multi-stage screening plan A, C and E maintain a relatively high score  
284 (the average value of EF1% for plan A is 7.80, for plan C is 7.68, for plan E is 7.55) in  
285 comparison with the exp\_z\_score method (the average value of EF1% is 7.74, the p-  
286 value for EF1% between plan E and exp\_z\_score is 0.420). With the reduction of the

287 consumed time for the multi-stage screening scheme, EF1% had a slight downward  
288 trend, but still can keep a higher value. In addition, the EF1% of multi-stage screening  
289 with using consensus scoring (for plan E is 7.55) is higher than that without consensus  
290 scoring (for plan F is 6.62) (Fig. 5b and Table 6). These indicated that the multi-stage  
291 screening using consensus score would remarkably improve the virtual screening  
292 results.

293 BEDROC is an important indicator of early enrichment capacity ( $\alpha=80.5$ ). We  
294 also observed that multi-stage screening plan A, C and E with consensus scoring (the  
295 average value of BEDROC( $\alpha=80.5$ ) for plan A is 0.203, for plan C is 0.197, for plan E  
296 is 0.189) exhibited a higher score of BEDROC relative to exp\_z\_score (the average  
297 value of BEDROC( $\alpha=80.5$ ) is 0.204, the p-value for BEDROC( $\alpha=80.5$ ) between plan  
298 E and exp\_z\_score is 0.032)(Fig. 5c), but a worse score in terms of AUC (the average  
299 value of AUC for plan A is 0.666, for plan C is 0.659, for plan E is 0.657, for  
300 exp\_z\_score is 0.686. The p-value for AUC between plan E and exp\_z\_score is  
301 0.0005) and logAUC (the average value of logAUC for plan A is 0.276, for plan C is  
302 0.268, for plan E is 0.262, for exp\_z\_score is 0.290. The p-value for AUC between  
303 plan E and exp\_z\_score is 0.0002) indicators. These results demonstrated that multi-  
304 stage screening combined with consensus score had a great advantage in early  
305 enrichment ability which is of great significance for drug discovery.

306 There was no significant difference between Plan E and exp\_z\_score in the term  
307 of EF1% values ( $p=0.420$ ), while the average calculation speed of each small  
308 molecule of Plan E (7.04 second/cpu/molecule) is nearly 9 times faster than  
309 exp\_z\_score (63.51 second/cpu/molecule). Therefore, multi-stage screening combined  
310 with consensus scoring can use the advantages of consensus scoring to synthesize  
311 each single scoring function to improve the ability of virtual screening effect, while  
312 significantly reducing the computational resources consumed in the whole process of  
313 consensus scoring. The multi-stage screening combined with consensus score has  
314 higher early enrichment ability and less computational resource consumption, which  
315 is of great significance for large-scale drug screening.

316

### 317 **Relationship between scoring score and hit rate**

318 In addition, we also analyzed four single scoring functions exp\_z\_score function  
319 of 51 protein targets from both DUD-E protein structures and AlphaFold protein  
320 structures, and counted their distributions and hit rates of small molecules.

321 To understand the relationship between the scoring and hit rate of small  
322 molecules in various scoring methods, we calculated the hit rate of active small  
323 molecules in the corresponding score segment of each scoring method. We found that  
324 the smaller the score of the autodock4, idock, and vinardo scoring functions, the more  
325 likely the small molecule is to combine with the protein. Moreover, as the score  
326 decreases, the hit rate tends to increase (Fig. 6a, b, e, f, i, j blue and green curves).  
327 While, according to the core algorithms of rf\_score and exp\_z\_score functions, the  
328 higher values demonstrated that these small molecules were more likely to combine  
329 with the protein. With the increase of the score, the hit rate tends to increase in terms  
330 of rf\_score and exp\_z\_score (Fig. 6c, d, g, h blue and green curves). For different  
331 methods, small molecule compounds with better scores are more likely to be active  
332 small molecule compounds and deserve more attention. This is also one of the bases  
333 for multi-stage screening to select the top small molecules for further screening in the  
334 next stage.

335 Since the number of molecules in some fraction segments of higher and lower  
336 scores is small, the hit rate error corresponding to the fraction segment may be large.  
337 So there is some fluctuation on the curve. The curve of score and hit rate is of great  
338 significance for evaluating the scores given by different methods. It is helpful for  
339 researchers to set threshold values based on scores and give priority to small  
340 molecules that are more likely to be active compounds. For example, if the score of a  
341 small molecule on a protein is less than -9 by using the autodock4 scoring function for  
342 docking, the small molecule is more likely to be an active small molecule.

343

### 344 **Discussion**

345 We used the targets in the DUD-E database, which is widely used in the  
346 benchmark test of virtual screening, and selected their corresponding AlphaFold  
347 structure and protein Apo structure. We applied a variety of docking software,  
348 multiple scoring functions and consensus scoring methods to perform virtual  
349 screening. To evaluate AlphaFold protein structure for virtual screening, we compared  
350 the performance of the protein structure of AlphaFold protein and DUD-E database,  
351 which represents the protein Holo structure.

352 We found that in the test of 51 selected targets, the virtual screening effect of  
353 protein structure in DUD-E database is better than that of AlphaFold protein structure  
354 as a whole, but there are 18 targets (about 35%) whose AlphaFold protein structure is  
355 equivalent to that of DUD-E protein, and 13 targets (about 25%) whose AlphaFold  
356 protein structure is better than that of DUD-E. In the test of 23 selected targets, it was  
357 found that the virtual screening effect of protein structure in DUD-E database was  
358 slightly better than that of AlphaFold protein structure, and AlphaFold protein  
359 structure was slightly better than that of Apo protein structure. Our results have  
360 guiding significance for the selection of protein structures in virtual screening.  
361 Generally, when the protein Holo structure is available, the protein Holo structure is  
362 preferred. When the protein Holo structure is not available, the AlphaFold protein  
363 structure can be used preferentially.

364 At present, the AlphaFold protein structure database has more than 200 million  
365 entries, providing human proteome and proteome of 47 other organisms. Although  
366 AlphaFold is only used to predict the monomer protein structure, it can be used to  
367 predict the protein composed of multiple monomers<sup>29</sup>. In addition, you can also use  
368 ColabFolder<sup>30</sup>, which is more user-friendly and easy to use, instead of installing the  
369 AlphaFold program and downloading the required database locally. Combined with  
370 our conclusions, the super large-scale protein structure database and low-cost and fast  
371 accurate method from sequence to three-dimensional protein structure provided by  
372 AlphaFold are of excellent use value for virtual screening, and of great significance  
373 for speeding up drug discovery and reducing drug development costs.

374 We tested 13 consensus scoring methods and found the `exp_z_score` was  
375 generally better than other methods, which showed to be higher than ECR in terms of  
376 BEDROC ( $\alpha=20.0$ ) and EF10%. The advantage of `exp_z_score` consensus method  
377 may be that it used a more reasonable Z-score standardized method and exponential  
378 function that was more regular for virtual screening. The Z-score standardization  
379 method can standardize the scoring data group of the scoring function for the small  
380 molecule set to the data group with a mean value of 0 and a standard deviation of 1.  
381 Compared with the standardized method of Rank, the Z-score standardization method  
382 can better retain the difference information of the same scoring function for different  
383 small molecules. In addition, the Z-score standardization method is less affected by  
384 the maximum and minimum values than the AASS standardization method.

385 With the rapid expansion of the small molecule compound library available for  
386 virtual screening and the super extensive protein structure database provided by  
387 AlphaFold, the large-scale virtual screening will play an increasingly important role in  
388 drug discovery. To reduce the computing time of virtual screening, improve the  
389 computing efficiency, save computing resources, and be able to use a variety of  
390 different scoring functions and consensus scoring methods, we propose a multi-stage  
391 screening combining consensus scoring methods. The average value of EF1% has no  
392 significant difference compared with consensus scoring, but its speed is about 9 times  
393 that of consensus scoring. Multi-stage screening combined with consensus score can  
394 flexibly adjust the number of small molecules screened at each stage to adjust the time  
395 and screening effect. It only uses a small amount of computing resources and has a  
396 high early enrichment capacity. Compared with the traditional multi-stage screening,  
397 the multi-stage screening combined with consensus score has great advantages in the  
398 early enrichment ability. Multi-stage screening combined with consensus score is used  
399 for large-scale drug screening, which is of great significance for speeding up drug  
400 discovery and reducing drug development costs.

401 We tried to optimize the AlphaFold structure by molecular dynamics in order to  
402 improve the virtual screening. This seems like a potential approach from the results of  
403 weel target shown in the Supplementary Information Table 20. But we can't split the

404 good structures from other structures for virtual screening without the active and  
405 decoy label. In other words, we can't split the good structures from other structures by  
406 some simple indication such as the best score or the average score. Our results and  
407 conclusions are similar to Nichols et al. In their studies they found that molecular  
408 dynamics can improve virtual screening results but they can't identify any evident  
409 relationship between characterization of molecular dynamics snapshots and virtual  
410 screening results<sup>31</sup>.

411 In the process of determining the structure of the protein used, we found that  
412 some proteins, such as HIVRT, had apparent binding pocket structures in their Holo  
413 and Apo structures around the positions of small co-crystalline molecules. But there  
414 was no obvious binding pocket structure in the AlphaFold protein structure. However,  
415 some proteins, such as AKT2, had obvious binding pocket structures in the Holo  
416 structure and AlphaFold structure around the eutectic small molecules, but there was  
417 no obvious binding pocket structure in the Apo structure. This indicated the typical  
418 inductive fit effect between proteins and small molecules. However, among the targets  
419 we examined, most AlphaFold protein structures and Apo protein structures have  
420 obvious binding pocket structures, which may be related to the flexibility of amino  
421 acid residues from protein binding pocket. These results also suggested the limitation  
422 of AlphaFold's prediction of protein structure: it was difficult to fully consider the  
423 impact of non-protein parts on protein structure, and it was unable to give non-protein  
424 structures such as water molecules, coenzymes, metal ions, etc. Since AlphaFold did  
425 not give a non-protein structure, we have removed the non-protein structure part from  
426 all tested protein structures.

427 We use AutoSite<sup>32</sup> to determine the search space, which has the advantage that  
428 we can flexibly describe the binding pocket structure of proteins. However, the  
429 multiple binding site structures given by AutoSite still need to be selected according  
430 to the position of small co-crystalline molecules relative to the protein structure. In  
431 addition, the combination pocket structure provided by AutoSite may be too large or  
432 too small. To avoid the influence of human operation on the experimental results, we

433 did not adjust the search space separately when the combined bag structure was too  
434 large or too small.

435 For our research, there are still some limitations. The number of targets used to  
436 compare protein Holo structure and AlphaFold structure is 51, and the number of  
437 targets used to compare protein Holo structure, protein Apo structure and AlphaFold  
438 structure is 23. Although some valuable data have been obtained from these targets,  
439 more targets can be used for testing to obtain more reliable data. In the treatment  
440 before the virtual screening of protein Holo structure, protein Apo structure and  
441 AlphaFold structure, we removed the non-protein structural parts from all structures.  
442 In the practical application of virtual screening, the non-protein components in protein  
443 Holo structure and protein Apo structure, such as water molecules, coenzymes, and  
444 metal ions, may be considered. The docking used in the study regards protein  
445 structure as rigid, while small molecules as flexible. Therefore, flexible docking<sup>33</sup>,  
446 that is, the amino acid residues in protein binding pockets are considered flexible and  
447 small molecules are considered flexible, which may be a method worth trying. The  
448 limitation of the curve drawn based on the relationship between scoring and hit rate  
449 lies in the high proportion of active small molecules in the DUD-E database and the  
450 small number of small molecules in the small molecule set, which is quite different  
451 from the large-scale virtual screening. In the test of consensus scoring method, we  
452 only considered the scoring of small molecules by the scoring function of docking  
453 software, but did not consider docking conformation. Taking into account docking  
454 conformation for consensus scoring may be a method worth trying. In the multi-stage  
455 screening test, the number of stages divided is three, and only three groups of schemes  
456 are set according to the proportion of small molecules entering the next stage in each  
457 stage. What's more, the small molecule set used for testing has fewer small molecules,  
458 which is several orders of magnitude less than the small molecule library of super  
459 large-scale virtual screening. Therefore, although our research can provide some  
460 guidance for super large-scale virtual screening, it is difficult to give a specific  
461 implementation plan for large-scale virtual screening. Further work is needed to  
462 develop a specific implementation plan for large-scale virtual screening.



463

## 464 **Methods**

### 465 **Targets**

466 To evaluate the effectiveness of AlphaFold protein structure in virtual screening,  
467 we used the targets in the DUD-E database. The DUD-E database contains 102  
468 protein targets, each of which includes protein structure (from the PDB database), co-  
469 crystalline small molecules to determine the docking software search space, active  
470 small molecules and decoy small molecules. It is widely utilized to assess the ability  
471 of virtual screening to distinguish active compounds from decoy compounds.

472 To ensure the diversity of protein types and the binding site of targets, we  
473 selected 51 representative targets in the DUD-E database as our test targets  
474 representing eight types of protein targets namely protease, nuclear receiver, kinase,  
475 cytochrome P450, ion channel, GPCR, other enzymes and miscellaneous proteins (Fig.  
476 1a). In total, more than 400,000 small molecules which interact with these 51 targets  
477 were used.

478 We tried to collect the Apo structure for the 51 selected targets in DUD-E  
479 database. However, only 23 out of these 51 targets have Apo structure available in the  
480 PDB database with known binding pocket. The information including the PDB ID and  
481 the resolution of 23 Apo structures can be find in the Supplementary Information  
482 Table 1.

483

### 484 **ColabFold**

485 Milot Mira et al. developed ColabFold based on Google Colab in order to enable  
486 researchers without relevant hardware resources to use AlphaFold2<sup>30</sup>.

487 Because HIVPR and HIVRT are non-monomer proteins, their protein structures  
488 cannot be directly obtained in the AlphaFold Database, so ColabFold is used for  
489 structure prediction. The input sequence and configuration can be find in the  
490 Supplementary Information.

491

### 492 **Protein structure preparation**

493        Considering the calculation time of virtual screening, the diversity of protein  
494 types and whether the AlphaFold protein structure can be used for virtual screening  
495 (structure availability and whether there is an obvious binding site), the DUD-E  
496 Database is selected( <http://dude.docking.org/targets> )to 51 targets in are used as  
497 virtual screening test targets.

498        The protein structure provided by DUD-E database corresponding to the 51  
499 selected targets is used as the test target of protein Holo structure. There are 50  
500 AlphaFold protein structures corresponding to 51 targets from AlphaFold  
501 Database( <https://alphafold.com/> ). HIVPR protein is predicted by ColabFold. The  
502 AlphaFold protein structure corresponding to all 51 targets serves as the virtual  
503 screening test target.

504        Because some proteins cannot or are difficult to find their Apo structures in the  
505 PDB database, and some proteins have no apparent binding sites, only 23 of the 51  
506 targets found available protein Apo structures as virtual screening test targets.

507        Each protein structure that is the target of the virtual screening test is processed  
508 as follows:

509        1. Use PDBFixer 1.8.1( <https://github.com/openmm/pdbfixer> ), replace non-  
510 standard residues, remove heterologous substances (including water molecules,  
511 coenzymes, metal ions, etc.), and add missing atoms to obtain the acceptor\_ fixed.pdb.

512        2. Use the script provided by AutoDock Vina to execute the following commands:  
513        prepare\_receptor -r receptor\_ fixed. pdb -o receptor. pdbqt -A hydrogens

514        The final receiver Pdbqt is the prepared protein file, which can be directly used  
515 as the input of AutoDock Vina, Qvina2 and idock to perform docking.

516

### 517 **Compounds preparation**

518        We get the active small molecule files corresponding to each target from DUD-E  
519 Database\_ final.sdf.gz and decoy small molecule file decoys\_ final.sdf.gz. Next,  
520 decompress and use openbabel 3.1.0<sup>34</sup> ( <https://github.com/openbabel/openbabel> ),  
521 perform split conversion to obtain the mol2 format file corresponding to each small  
522 molecule, and then use the script provided by AutoDock Vina for each small molecule

523 to execute the following commands for hydrogenation and preparation of small  
524 molecule file formats:

525 `prepare_ligand -l molecule. mol2 -o molecule. pdbqt -A hydrogens`

526 The final obtained Molecule Pdbqt is the prepared small molecule file, which can  
527 be directly used as the input of AutoDock Vina, Qvina2 and idock to perform docking.

528

### 529 **Search space**

530 AutoSite 1.0.0<sup>32</sup> is used to predict binding sites for each protein structure that is  
531 the target of the virtual screening test, and PyMOL (TM) Molecular Graphics System,  
532 Version 2.6.0a0 Open-Source<sup>35</sup> is used to visually inspect and select binding sites.

533 For the protein structure in the DUD-E database, refer to the eutectic small  
534 molecule crystal in the DUD-E database\_ The position of ligand.mol2 selects the  
535 binding site.

536 For the AlphaFold protein structure and Apo structure, align the protein structure  
537 and small co-crystalline molecules corresponding to the target in the DUD-E database  
538 with the AlphaFold protein structure and Apo structure. Then, select the binding site  
539 by referring to the position of small co-crystalline molecules.

540 After determining the protein structure binding site, use the PandasPdb module  
541 of biopandas 0.2.9 to calculate the spatial coordinate information of the binding site  
542 structure, including size and location, and then add 8 Å to each of the x, y, and z  
543 dimensions of the binding site structure size as the search space for AutoDock Vina,  
544 Qvina2, and idock.

545

### 546 **Docking**

547 Three docking software AutoDock Vina 1.2.3, QuickVina 2.1 (Qvina2) and idock  
548 2.2.3 are used. And five scoring functions including autodock4 and vinarado in  
549 AutoDock Vina, qvina scoring function, idock scoring function, and rf\_scoring  
550 function in idock were employed. For each docking, the seed used is 20011204, and  
551 only the best docking position is considered as output. The exhaustiveness value used

552 by AutoDock Vina and Qvina2 is 1, and the tasks used by idock are 8. Neither flexible  
553 docking nor hydration docking is used.

554 The used autodock4 (ad4) scoring function and vinardo scoring function are the  
555 scoring functions provided by AutoDock Vina 1.2.3. When using the autodock4  
556 scoring function and vinardo scoring function, perform docking respectively to obtain  
557 the scoring of a single small molecule by autodock4 scoring function and the scoring  
558 of a single small molecule by vinardo scoring function. Affinity map is also required  
559 when using the autodock4 scoring function. Each docking is prepared by autogrid4  
560 (AutoGrid 4.2.7. x.2019-07-11) and AutoDock Vina\_ The affinity graph was  
561 calculated by gpf.py.

562 The qvina scoring function used is QuickVina 2.1 scoring function.

563 The score function used by idock and rf\_score is provided by idock 2.2.3. When  
564 idock 2.2.3 is executed, only one time of docking would be performed to obtain the  
565 score of idock and rf\_score functions for a single small molecule screening.

566

### 567 **Consensus scoring**

568 For every target, to integrate different scoring functions, different scoring  
569 functions may need to be standardized. Common standardization methods include  
570 Rank, AASS(Average of auto-scaled scores), Z-score, etc. Only after the scores  
571 between different scoring functions are standardized can they be reasonably  
572 calculated and compared.

573 According to the direction of scoring, the scoring function can be divided into  
574 two categories. The first category is that the larger the scoring value represents the  
575 stronger the binding ability of small molecule compounds to protein. The second  
576 category is that the higher the score represents the weaker the binding ability of small  
577 molecule compounds to proteins is.

578 The effective scores of the first category are all scores with a score value greater  
579 than 0; The second type of valid score is all scores with a score value less than 0.

### 580 **Rank**

581 Rank as the standardized score.

582 For the first category, the scores are sorted from the largest to the smallest.

583 For the second category, the scores are sorted from the smallest to the largest.

584 That is, the scores from two categories are ranked from good to bad. And then, use

585 rank as standardized score for the two categories. The smaller the standardized score,

586 the better.

### 587 **Z-score**

588 Standardize score  $S_{ij}$  for small molecule  $j$  on scoring function  $i$  by Z-score to obtain

589 the standardized score  $Z_{ij}$ .

590 The mean of the effective scores on the scoring function  $i$  is  $\bar{S}_i$ , and the standard

591 deviation of that is  $\sigma_i$ .

592 The scores from the first category are standardized by the following formula:

$$Z_{ij} = \frac{S_{ij} - \bar{S}_i}{\sigma_i}$$

593 The scores from the second category are standardized by the following formula:

$$Z_{ij} = \frac{\bar{S}_i - S_{ij}}{\sigma_i}$$

594 For the two categories, the larger the standardized score, the better.

### 595 **AASS**

596 Standardize score  $S_{ij}$  for small molecule  $j$  on scoring function  $i$  by AASS to obtain

597 the standardized score  $A_{ij}$ .

598 For the first category and the scoring function  $i$ , the largest score is  $Best_i$ , the

599 smallest score is  $Worst_i$ .

600 For the second category and the scoring function  $i$ , the smallest score is  $Best_i$ , the

601 largest score is  $Worst_i$ .

602 That is, for the two categories, the best score is  $Best_i$ , the worst is  $Worst_i$ .

603 The scores from the two categories are standardized by the following formula:

$$A_{ij} = \frac{S_{ij} - Worst_i}{Best_i - Worst_i}$$

604 For the two categories, the larger the standardized score, the better.

605 At the same time, another key to consensus scoring is how to process standardized  
606 scores to obtain the final consensus score.

607 We selected four calculation methods: Summation (Sum), Best, Worst, and  
608 Summation after Exponentiation (Exp).

609 **Sum**

610 The standardized score for small molecule  $j$  on scoring function  $i$  is  $V_{ij}$ . The  
611 consensus score  $Sum_j$  for small molecule  $j$  is calculated by the following formula:

$$Sum_j = \sum_i V_{ij}$$

612 **Best**

613 The consensus score  $Best_j$  for small molecule  $j$  is the best standardized score in every  
614 scoring function.

615 **Worst**

616 The consensus score  $Worst_j$  for small molecule  $j$  is the worst standardized score in  
617 every scoring function.

618 **Exp**

619 The standardized score for small molecule  $j$  on scoring function  $i$  is  $V_{ij}$ . The  
620 consensus score  $Exp_j$  for small molecule  $j$  is calculated by the following formula:

$$Exp_j = \sum_i e^{V_{ij}}$$

621 **Consensus function**

622 According to the combination of standardization methods (3) and calculation methods  
623 (4), 12 consensus scoring methods were obtained. The combination and name are as  
624 follows:

	Rank	Z-score	AASS
Sum	sum_rank	sum_z_score	sum_aass
Best	best_rank	best_z_score	best_aass
Worst	worst_rank	worst_z_score	worst_aass
Exp	exp_rank	exp_z_score	exp_aass

625 For example, the complete description of `exp_z_score` is as follows:

626 The mean of the effective scores on the scoring function  $i$  is  $\bar{S}_i$ , and the standard  
627 deviation of that is  $\sigma_i$ . The standardized score for small molecule  $j$  on scoring  
628 function  $i$  is  $Z_{ij}$ .

629 The scores from the first category are standardized by the following formula:

$$Z_{ij} = \frac{S_{ij} - \bar{S}_i}{\sigma_i}$$

630 The scores from the second category are standardized by the following formula:

$$Z_{ij} = \frac{\bar{S}_i - S_{ij}}{\sigma_i}$$

631 The final consensus score  $Exp_j$  for small molecule  $j$  is calculated by the following  
632 formula:

$$Exp_j = \sum_i e^{Z_{ij}}$$

633 The larger the final consensus score, the more likely it is an active small molecule.

634 While `exp_rank(ECR)` has changed the method of `Exp`:The rank for small molecule  $j$   
635 on scoring function  $i$  is  $R_{ij}$ . The number of the effective scores on the scoring  
636 function  $i$  is  $n_i$ . The consensus score  $Exp_j$  for small molecule  $j$  is calculated by the  
637 following formula:

$$Exp_j = \sum_i e^{-R_{ij}/(n_i*10\%)}$$

638 The larger the final consensus score, the more likely it is an active small molecule.

639 In addition, Rank by Vote (RBV) consensus scoring method was tested.

640

#### 641 **RBV**

642 The threshold value is set to  $x=\text{top}10\%$  and the ranking of small molecule  $j$  in the  
643 scoring function  $i$  is within the threshold value  $x$ , one vote will be obtained. The sum  
644 of votes of small molecule  $j$  on all scoring functions is the final score. Small  
645 molecules with the same number of votes are randomly ranked.

646

647 **Multi-stage virtual screening with consensus ranking**

648 Prepare small molecules, proteins and search space, and then start multi-stage  
649 screening. In the first stage, idock docking is used for virtual screening to obtain idock  
650 rank and rf\_score rank, and then ir rank is obtained through the consensus score of  
651 idock rank and rf\_score rank. The top x% of the ir rank are taken to enter the second  
652 stage. In the second stage, we use the vinardo scoring function to conduct virtual  
653 screening to obtain the vinardo rank, and then take the small molecules in the topx%  
654 of the ir rank in the idock rank and rf\_score rank and all the small molecules in the  
655 vinardo rank for consensus scoring to obtain the irv rank. The top y% of the irv rank  
656 are taken to enter the third stage. In the third stage, the autodock4 scoring function is  
657 used to conduct virtual screening, and the autodock4 rank is obtained. Then the small  
658 molecules in idock rank, rf\_score rank and vinardo rank that are y% before irv rank  
659 and all small molecules in autodock4 rank are scored by consensus to obtain irva rank.  
660 The final ranking is obtained by ranking the small molecules in irva rank, the small  
661 molecules in irv rank that are not in irva rank, and the small molecules in irv rank that  
662 are not in irv rank according to the original order within the ranking. The x and y  
663 parameters can be adjusted flexibly. Three groups of schemes are set, Plan A(x=40 ,  
664 y=50) , Plan C(x=20 , y=50) , Plan E(x=10 , y=50).

665

666 **Multi-stage virtual screening without consensus ranking**

667 Prepare small molecules, proteins and search space, and then start multi-stage  
668 screening. In the first stage, idock docking is used for virtual screening to obtain idock  
669 rank, the top x% of the idock rank are taken to enter the second stage. In the second  
670 stage, we use the vinardo scoring function to conduct virtual screening to obtain the  
671 vinardo rank, the top y% of the vinardo rank are taken to enter the third stage. In the  
672 third stage, the autodock4 scoring function is used to conduct virtual screening, and  
673 the autodock4 rank is obtained. The final ranking is obtained by ranking the small  
674 molecules in autodock4 rank, the small molecules in vinardo rank that are not in  
675 autodock4 rank, and the small molecules in idock rank that are not in vinardo rank  
676 according to the original order within the ranking. The x and y parameters can be



677 adjusted flexibly. Three groups of schemes are set, Plan B(x=40 , y=50) , Plan  
678 D(x=20 , y=50) , Plan F(x=10 , y=50).

679

## 680 **Metrics**

681 We used AUC (area under the receiver operating characteristic curve), logAUC  
682 (area under the semilogarithmic receiver operating characteristic curve), EF1%  
683 (enrichment factor 1%), EF5%, EF10% and BEDROC (Boltzmann-enhanced  
684 discrimination of receiver operating characteristic,  $\alpha= 321.9$ ,  $\alpha= 80.5$ ,  $\alpha= 20.0$ ) as  
685 benchmark metrics to assess the result. LogAUC is the area under the semi-  
686 logarithmic ROC curve, and the early enrichment capacity has a greater weight in  
687 logAUC than AUC<sup>12</sup>. Enrichment factors (1%, 5%, and 10%) are used to truncate the  
688 assessment of early enrichment capacity under specific thresholds, which are closely  
689 related to the practical application of virtual screening. The BEDROC<sup>36</sup> was proposed  
690 to consider all small molecules by adjusting  $\alpha$  value for the weight of early  
691 enrichment capacity. The greater the early enrichment capacity, the larger  $\alpha$  weight  
692 value.

693

## 694 **Hits rate**

695 In a collection of small molecules, the number of active compounds in the collection  
696 is *actives*, the number of all compounds in the collection is *total*. *Hits rate* is  
697 calculated by the following formula:

$$Hits\ rate = \frac{actives}{total}$$

698 For a collection of small molecules in the top x% of a rank, the number of active  
699 compounds in the collection is *actives(x%)*, the number of all compounds in the  
700 collection is *total(x%)*. *Hits rate(x%)* is calculated by the following formula:

$$Hits\ rate(x\%) = \frac{actives(x\%)}{total(x\%)}$$

## 701 **Enrichment factor**

702 For a collection of small molecules in the top x% of a rank, the enrichment factor  
703 x%(EFx%) is calculated by the following formula:

$$EFx\% = \frac{Hits\ rate(x\%)}{Hits\ rate(100\%)}$$

704 **BEDROC**

705 BEDROC is defined as follows:

$$BEDROC(\alpha) = \frac{\sum_{i=1}^n e^{-\alpha r_i/N}}{\frac{n}{N} \left( \frac{1 - e^{-\alpha}}{e^{\alpha/N} - 1} \right)} \times \frac{R_a \sinh(\alpha/2)}{\cosh(\alpha/2) - \cosh(\alpha/2 - \alpha R_a)} + \frac{1}{1 - e^{\alpha(1-R_a)}}$$

706 While the  $n$  is the number of active compounds, the  $N$  is the number of all  
707 compounds, the rate  $R_a = n/N$ , the  $r_i$  is the rank of  $i$  th active compound<sup>36</sup>.

708 **Receiver Operating Characteristics (ROC) Curve and its area under the curve**  
709 **(AUC)**

710 The Y-axis is the true positive rate (TPR, also named actives found rate) under a  
711 specific threshold. The X-axis is the false positive rate (FPR, also named decoys  
712 found rate) under a specific threshold.

713 TPR is calculated by the following formula:

$$TPR = \frac{TP}{TP + FN}$$

714 FPR is calculated by the following formula:

$$FPR = \frac{FP}{FP + TN}$$

715 At a specific threshold, TP is the number of true positive compounds, FN is the  
716 number of false negative compounds, FP is the number of false positive compounds,  
717 and TN is the number of true negative compounds.

718 The area under the ROC curve is AUC, and its value range is [0,1]. The AUC  
719 corresponding to random screening is 0.5.

720 **Semilogarithmic Receiver Operating Characteristics (ROC) Curve and its area**  
721 **under the curve (logAUC)**

722 After drawing ROC curve, let  $\lambda = 0.001$ . Points with abscissa less than  $\lambda$  are ignored  
723 and the abscissa of other points is performed logarithmic( $\log_{10}$ ) operation. That is,  
724 the point  $(x_i, y_i)$  in the ROC curve, the point  $(\log_{10}(x_i), y_i)$  in the semilogarithmic  
725 ROC.

726 The abscissa value range of Semilogarithmic ROC curve is [-3,0]. The area of  
727 semilogarithmic ROC curve is  $area$ . Let  $\log AUC = area/3$ , and its value range is  
728 [0,1].

729 **Time**

730 Use the Linux time command to count the time, and use user time+sys time as the  
731 time used for virtual screening.

### 732 **System**

733 All docking calculations are performed in containers created by images built with  
734 Docker 20.10.17 software. Docker is an open source lightweight virtualization  
735 technology (<https://www.docker.com/>).

736 The host system is Linux 5.4.0-131 generic.

### 737 **CPU**

738 Intel(R) Xeon(R) Gold 6238R CPU @ 2.20GHz.

### 739 **Figure**

740 Violin chart, ROC curve, scatter plot, fitting line and histogram are drawn with  
741 matplotlib 3.6.1<sup>37</sup>.

742 The three-dimensional structure diagram of small molecules and protein is drawn with  
743 PyMOL (TM) Molecular Graphics System, Version 2.6.0a0 Open-Source<sup>35</sup>.

744

### 745 **Funding**

746 This work was supported by National Natural Science Foundation of China  
747 (82260488 and 32200679), China Postdoctoral Science Foundation (2021TQ0137,  
748 2021M701544), Yunnan Applied Basic Research Key Projects (202001AT070102,  
749 and 202001AT070104), Natural Science Foundation of Chongqing  
750 (CSTB2022NSCQ-MSX056), Postgraduate innovation special fund project of Jiangxi  
751 (YC2022—s047), College Students' Innovative Entrepreneurial Training Plan  
752 Program of Nanchang University (2022CX024, 2022CX158), Scientific Research  
753 Training Program of Nanchang University (2022).

754

### 755 **Conflict of interest**

756 The authors declare that they have no conflict of interest.

757

### 758 **Reference**

759 (1) Jumper, J.; Evans, R.; Pritzel, A.; Green, T.; Figurnov, M.; Ronneberger, O.;

- 760 Tunyasuvunakool, K.; Bates, R.; Žídek, A.; Potapenko, A.; Bridgland, A.; Meyer,  
761 C.; Kohl, S. A. A.; Ballard, A. J.; Cowie, A.; Romera-Paredes, B.; Nikolov, S.;  
762 Jain, R.; Adler, J.; Back, T.; Petersen, S.; Reiman, D.; Clancy, E.; Zielinski, M.;  
763 Steinegger, M.; Pacholska, M.; Berghammer, T.; Bodenstein, S.; Silver, D.;  
764 Vinyals, O.; Senior, A. W.; Kavukcuoglu, K.; Kohli, P.; Hassabis, D. Highly  
765 Accurate Protein Structure Prediction with AlphaFold. *Nature* **2021**, *596* (7873),  
766 583–589. <https://doi.org/10.1038/s41586-021-03819-2>.
- 767 (2) Binder, J. L.; Berendzen, J.; Stevens, A. O.; He, Y.; Wang, J.; Dokholyan, N. V.;  
768 Oprea, T. I. AlphaFold Illuminates Half of the Dark Human Proteins. *Current*  
769 *Opinion in Structural Biology* **2022**, *74*, 102372.  
770 <https://doi.org/10.1016/j.sbi.2022.102372>.
- 771 (3) Ros-Lucas, A.; Martinez-Peinado, N.; Bastida, J.; Gascón, J.; Alonso-Padilla, J.  
772 The Use of AlphaFold for In Silico Exploration of Drug Targets in the Parasite  
773 *Trypanosoma Cruzi*. *Frontiers in Cellular and Infection Microbiology* **2022**, *12*.
- 774 (4) Weng, Y.; Pan, C.; Shen, Z.; Chen, S.; Xu, L.; Dong, X.; Chen, J. Identification of  
775 Potential WSB1 Inhibitors by AlphaFold Modeling, Virtual Screening, and  
776 Molecular Dynamics Simulation Studies. *Evidence-Based Complementary and*  
777 *Alternative Medicine* **2022**, *2022*, e4629392.  
778 <https://doi.org/10.1155/2022/4629392>.
- 779 (5) Yang, C.; Alam, A.; Alhumaydhi, F. A.; Khan, M. S.; Alsagaby, S. A.; Al  
780 Abdulmonem, W.; Hassan, M. I.; Shamsi, A.; Bano, B.; Yadav, D. K. Bioactive  
781 Phytoconstituents as Potent Inhibitors of Tyrosine-Protein Kinase Yes (YES1):  
782 Implications in Anticancer Therapeutics. *Molecules* **2022**, *27* (10), 3060.  
783 <https://doi.org/10.3390/molecules27103060>.
- 784 (6) Wong, F.; Krishnan, A.; Zheng, E. J.; Stärk, H.; Manson, A. L.; Earl, A. M.;  
785 Jaakkola, T.; Collins, J. J. Benchmarking AlphaFold-Enabled Molecular Docking  
786 Predictions for Antibiotic Discovery. *Molecular Systems Biology* **2022**, *18* (9),  
787 e11081. <https://doi.org/10.15252/msb.202211081>.
- 788 (7) Scardino, V.; Filippo, J. I. D.; Cavasotto, C. How Good Are AlphaFold Models for  
789 Docking-Based Virtual Screening? **2022**. [https://doi.org/10.26434/chemrxiv-](https://doi.org/10.26434/chemrxiv-2022-sgj8c)  
790 [2022-sgj8c](https://doi.org/10.26434/chemrxiv-2022-sgj8c).
- 791 (8) Zhang, Y.; Vass, M.; Shi, D.; Abualrous, E.; Chambers, J.; Chopra, N.; Higgs, C.;  
792 Kasavajhala, K.; Li, H.; Nandekar, P.; Sato, H.; Miller, E.; Repasky, M.; Jerome, S.  
793 Benchmarking Refined and Unrefined AlphaFold2 Structures for Hit Discovery.  
794 **2022**. <https://doi.org/10.26434/chemrxiv-2022-kcn0d-v2>.
- 795 (9) Alon, A.; Lyu, J.; Braz, J. M.; Tummino, T. A.; Craik, V.; O’Meara, M. J.; Webb,  
796 C. M.; Radchenko, D. S.; Moroz, Y. S.; Huang, X.-P.; Liu, Y.; Roth, B. L.; Irwin, J.  
797 J.; Basbaum, A. I.; Shoichet, B. K.; Kruse, A. C. Structures of the  $\Sigma 2$  Receptor  
798 Enable Docking for Bioactive Ligand Discovery. *Nature* **2021**, *600* (7890), 759–  
799 764. <https://doi.org/10.1038/s41586-021-04175-x>.
- 800 (10) Gorgulla, C.; Boeszoermyeni, A.; Wang, Z.-F.; Fischer, P. D.; Coote, P. W.;  
801 Padmanabha Das, K. M.; Malets, Y. S.; Radchenko, D. S.; Moroz, Y. S.; Scott, D.  
802 A.; Fackeldey, K.; Hoffmann, M.; Iavniuk, I.; Wagner, G.; Arthanari, H. An Open-  
803 Source Drug Discovery Platform Enables Ultra-Large Virtual Screens. *Nature*

- 804       **2020**, *580* (7805), 663–668. <https://doi.org/10.1038/s41586-020-2117-z>.
- 805 (11)Lyu, J.; Wang, S.; Balius, T. E.; Singh, I.; Levit, A.; Moroz, Y. S.; O’Meara, M. J.;  
806       Che, T.; Alga, E.; Tolmacheva, K.; Tolmachev, A. A.; Shoichet, B. K.; Roth, B.  
807       L.; Irwin, J. J. Ultra-Large Library Docking for Discovering New Chemotypes.  
808       *Nature* **2019**, *566* (7743), 224–229. <https://doi.org/10.1038/s41586-019-0917-9>.
- 809 (12)Bender, B. J.; Gahbauer, S.; Luttens, A.; Lyu, J.; Webb, C. M.; Stein, R. M.; Fink,  
810       E. A.; Balius, T. E.; Carlsson, J.; Irwin, J. J.; Shoichet, B. K. A Practical Guide to  
811       Large-Scale Docking. *Nat Protoc* **2021**, *16* (10), 4799–4832.  
812       <https://doi.org/10.1038/s41596-021-00597-z>.
- 813 (13)Trott, O.; Olson, A. J. AutoDock Vina: Improving the Speed and Accuracy of  
814       Docking with a New Scoring Function, Efficient Optimization, and  
815       Multithreading. *Journal of Computational Chemistry* **2010**, *31* (2), 455–461.  
816       <https://doi.org/10.1002/jcc.21334>.
- 817 (14)Eberhardt, J.; Santos-Martins, D.; Tillack, A. F.; Forli, S. AutoDock Vina 1.2.0:  
818       New Docking Methods, Expanded Force Field, and Python Bindings. *J. Chem. Inf.*  
819       *Model.* **2021**, *61* (8), 3891–3898. <https://doi.org/10.1021/acs.jcim.1c00203>.
- 820 (15)Alhossary, A.; Handoko, S. D.; Mu, Y.; Kwok, C.-K. Fast, Accurate, and Reliable  
821       Molecular Docking with QuickVina 2. *Bioinformatics* **2015**, *31* (13), 2214–2216.  
822       <https://doi.org/10.1093/bioinformatics/btv082>.
- 823 (16)Li, H.; Leung, K.-S.; Wong, M.-H. Idock: A Multithreaded Virtual Screening Tool  
824       for Flexible Ligand Docking. In *2012 IEEE Symposium on Computational*  
825       *Intelligence in Bioinformatics and Computational Biology (CIBCB)*; 2012; pp 77–  
826       84. <https://doi.org/10.1109/CIBCB.2012.6217214>.
- 827 (17)Abagyan, R.; Totrov, M.; Kuznetsov, D. ICM—A New Method for Protein  
828       Modeling and Design: Applications to Docking and Structure Prediction from the  
829       Distorted Native Conformation. *Journal of Computational Chemistry* **1994**, *15* (5),  
830       488–506. <https://doi.org/10.1002/jcc.540150503>.
- 831 (18)Halgren, T. A.; Murphy, R. B.; Friesner, R. A.; Beard, H. S.; Frye, L. L.; Pollard,  
832       W. T.; Banks, J. L. Glide: □ A New Approach for Rapid, Accurate Docking and  
833       Scoring. 2. Enrichment Factors in Database Screening. *J. Med. Chem.* **2004**, *47*  
834       (7), 1750–1759. <https://doi.org/10.1021/jm030644s>.
- 835 (19)Jones, G.; Willett, P.; Glen, R. C.; Leach, A. R.; Taylor, R. Development and  
836       Validation of a Genetic Algorithm for Flexible Docking<sup>11</sup>Edited by F. E. Cohen.  
837       *Journal of Molecular Biology* **1997**, *267* (3), 727–748.  
838       <https://doi.org/10.1006/jmbi.1996.0897>.
- 839 (20)Palacio-Rodríguez, K.; Lans, I.; Cavasotto, C. N.; Cossio, P. Exponential  
840       Consensus Ranking Improves the Outcome in Docking and Receptor Ensemble  
841       Docking. *Sci Rep* **2019**, *9* (1), 5142. <https://doi.org/10.1038/s41598-019-41594-3>.
- 842 (21)Preto, J.; Gentile, F. Assessing and Improving the Performance of Consensus  
843       Docking Strategies Using the DockBox Package. *J Comput Aided Mol Des* **2019**,  
844       *33* (9), 817–829. <https://doi.org/10.1007/s10822-019-00227-7>.
- 845 (22)Tuccinardi, T.; Poli, G.; Romboli, V.; Giordano, A.; Martinelli, A. Extensive  
846       Consensus Docking Evaluation for Ligand Pose Prediction and Virtual Screening  
847       Studies. *J. Chem. Inf. Model.* **2014**, *54* (10), 2980–2986.

- 848 <https://doi.org/10.1021/ci500424n>.
- 849 (23)Ochoa, R.; Palacio-Rodriguez, K.; Clemente, C. M.; Adler, N. S. DockECR:  
850 Open Consensus Docking and Ranking Protocol for Virtual Screening of Small  
851 Molecules. *Journal of Molecular Graphics and Modelling* **2021**, *109*, 108023.  
852 <https://doi.org/10.1016/j.jmgm.2021.108023>.
- 853 (24)Scardino, V.; Bollini, M.; Cavasotto, C. N. Combination of Pose and Rank  
854 Consensus in Docking-Based Virtual Screening: The Best of Both Worlds. *RSC*  
855 *Adv.* **2021**, *11* (56), 35383–35391. <https://doi.org/10.1039/D1RA05785E>.
- 856 (25)Sterling, T.; Irwin, J. J. ZINC 15 – Ligand Discovery for Everyone. *J. Chem. Inf.*  
857 *Model.* **2015**, *55* (11), 2324–2337. <https://doi.org/10.1021/acs.jcim.5b00559>.
- 858 (26)McGovern, S. L.; Shoichet, B. K. Information Decay in Molecular Docking  
859 Screens against Holo, Apo, and Modeled Conformations of Enzymes. *J. Med.*  
860 *Chem.* **2003**, *46* (14), 2895–2907. <https://doi.org/10.1021/jm0300330>.
- 861 (27)Lee, H. S.; Lee, C. S.; Kim, J. S.; Kim, D. H.; Choe, H. Improving Virtual  
862 Screening Performance against Conformational Variations of Receptors by Shape  
863 Matching with Ligand Binding Pocket. *J. Chem. Inf. Model.* **2009**, *49* (11), 2419–  
864 2428. <https://doi.org/10.1021/ci9002365>.
- 865 (28)Mysinger, M. M.; Carchia, M.; Irwin, John. J.; Shoichet, B. K. Directory of  
866 Useful Decoys, Enhanced (DUD-E): Better Ligands and Decoys for Better  
867 Benchmarking. *J. Med. Chem.* **2012**, *55* (14), 6582–6594.  
868 <https://doi.org/10.1021/jm300687e>.
- 869 (29)Evans, R.; O'Neill, M.; Pritzel, A.; Antropova, N.; Senior, A.; Green, T.; Žídek,  
870 A.; Bates, R.; Blackwell, S.; Yim, J.; Ronneberger, O.; Bodenstern, S.; Zielinski,  
871 M.; Bridgland, A.; Potapenko, A.; Cowie, A.; Tunyasuvunakool, K.; Jain, R.;  
872 Clancy, E.; Kohli, P.; Jumper, J.; Hassabis, D. Protein Complex Prediction with  
873 AlphaFold-Multimer. bioRxiv March 10, 2022, p 2021.10.04.463034.  
874 <https://doi.org/10.1101/2021.10.04.463034>.
- 875 (30)Mirdita, M.; Schütze, K.; Moriwaki, Y.; Heo, L.; Ovchinnikov, S.; Steinegger, M.  
876 ColabFold: Making Protein Folding Accessible to All. *Nat Methods* **2022**, *19* (6),  
877 679–682. <https://doi.org/10.1038/s41592-022-01488-1>.
- 878 (31)Nichols, S. E.; Baron, R.; Ivetac, A.; McCammon, J. A. Predictive Power of  
879 Molecular Dynamics Receptor Structures in Virtual Screening. *J. Chem. Inf.*  
880 *Model.* **2011**, *51* (6), 1439–1446. <https://doi.org/10.1021/ci200117n>.
- 881 (32)Ravindranath, P. A.; Sanner, M. F. AutoSite: An Automated Approach for Pseudo-  
882 Ligands Prediction—from Ligand-Binding Sites Identification to Predicting Key  
883 Ligand Atoms. *Bioinformatics* **2016**, *32* (20), 3142–3149.  
884 <https://doi.org/10.1093/bioinformatics/btw367>.
- 885 (33)Spyrakis, F.; Cavasotto, C. N. Open Challenges in Structure-Based Virtual  
886 Screening: Receptor Modeling, Target Flexibility Consideration and Active Site  
887 Water Molecules Description. *Archives of Biochemistry and Biophysics* **2015**, *583*,  
888 105–119. <https://doi.org/10.1016/j.abb.2015.08.002>.
- 889 (34)O'Boyle, N. M.; Banck, M.; James, C. A.; Morley, C.; Vandermeersch, T.;  
890 Hutchison, G. R. Open Babel: An Open Chemical Toolbox. *Journal of*  
891 *Cheminformatics* **2011**, *3* (1), 33. <https://doi.org/10.1186/1758-2946-3-33>.

- 892 (35) Schrödinger, LLC. The PyMOL Molecular Graphics System, Version 2.6.0a0  
893 Open-Source, 2021.
- 894 (36) Truchon, J.-F.; Bayly, C. I. Evaluating Virtual Screening Methods: □ Good and  
895 Bad Metrics for the “Early Recognition” Problem. *J. Chem. Inf. Model.* **2007**, *47*  
896 (2), 488–508. <https://doi.org/10.1021/ci600426e>.
- 897 (37) Hunter, J. D. Matplotlib: A 2D Graphics Environment. *Computing in Science &*  
898 *Engineering* **2007**, *9* (3), 90–95. <https://doi.org/10.1109/MCSE.2007.55>.

900

901

902

### 903 **Figure Legends**

904 **Figure 1. Virtual screening of protein structure in DUD-E database and**  
905 **AlphaFold.** (a) Number distribution of protein categories of 51 targets selected from  
906 the DUD-E database. (b) Draw a semi-logarithmic ROC curve based on the virtual  
907 screening results of 51 targets by the AutoDock4 scoring function ( $\lambda=0.001$ ). The red  
908 curve depicts the screening result of protein structure in the DUD-E database, and the  
909 blue curve depicts the AlphaFold protein structure screening result. (c) 51 targets in  
910 autodock4, idock, qvina, rf\_score, violin chart of screening results under the vinarado  
911 scoring function, using EF1%, AUC, logAUC, BEDROC ( $\alpha=80.5$ ) indicators are  
912 used to evaluate the screening effect, among which red is the screening result of  
913 protein structure in DUD-E database, and blue is the screening result of AlphaFold  
914 protein structure. According to paired t-test, two groups of data (red and blue) with the  
915 same scoring function and the same indicator have been marked with p values that are  
916 statistically different. \*\* represents  $p<0.01$ , and \* represents  $p<0.05$ .

917

918 **Figure 2. Binding pocket structure and virtual screening results of Holo,**  
919 **AlphaFold and Apo protein structures.** (a-b) Water molecules have been removed  
920 from protein structures. AlphaFold structure and Apo (obtained from PDB database)  
921 structure are aligned with Holo (obtained from DUD-E database) structure, displayed  
922 in grid form, and cut the same section. Green is Holo structure, purple is AlphaFold

923 structure, and yellow is Apo structure. PyMOL is used for operation and drawing.  
924 HIVRT (a) and AKT2 (b) protein structures of Holo, AlphaFold and Apo three-  
925 dimensional grid diagram sections were displayed. (c) 23 targets in autodock4, idock,  
926 rf\_score, violin chart of screening results under the vinardo scoring function, using  
927 EF1%, AUC, logAUC, BEDROC ( $\alpha= 80.5$ ) indicators are used to evaluate the  
928 screening effect. Among them, red group represent the screening result of protein  
929 structure (Holo) in DUD-E database, blue group represent screening result of  
930 AlphaFold protein structure, and green is the screening result of Apo protein structure.  
931 Three groups of data (red, blue and green) with the same scoring function and the  
932 same indicator are all marked according to paired t-test, and those with statistically  
933 different. \*\* represents  $p<0.01$ , and \* represents  $p<0.05$ .

934

935 **Figure 3. Evaluation of consensus scoring method based on score or ranking.** (a)

936 According to the DUD-E database, 51 targets are uploaded into autodock4, idock,  
937 qvina, rf\_Score and vinardo are the screening results of five scoring functions.  
938 Scatter plots and fitting lines are drawn between different scoring functions. The red  
939 is logAUC, and the blue is BEDROC ( $\alpha= 80.5$ ). (b) According to DUD-E database,  
940 BACE1 target is used as an example to draw scatter plots of four consensus scoring  
941 methods: exp\_z\_score, sum\_rank, best\_rank, and worst\_rank (all with idock and  
942 rf\_score as inputs). The abscissa is the ranking of small molecules in idock, and the  
943 ordinate is the ranking of small molecules in rf\_score, in which that gray and silver  
944 are decoys, red and purple are actives, and silver and purple are the top 10% of the  
945 ranking obtained by the corresponding consensus scoring method. (c) The 51 targets  
946 in the DUD-E database screen the violin chart of the results under 4 single scoring  
947 functions (blue) and 13 consensus scoring methods (red, with autodock4, idock,  
948 rf\_score, vinardo as the input), and use EF1%, AUC, logAUC, BEDROC ( $\alpha= 80.5$ )  
949 indicators are used to evaluate the screening effect. The x-axis scoring function is  
950 sorted by the average value of logAUC, with marked exp\_z\_score corresponding to  
951 the upper limit and lower limit of the data (gray dotted line), quartile (green dotted  
952 line), and median (blue dotted line).



953

954 **Figure 4. Re-evaluation of Holo, AlphaFold and Apo protein structures using**

955 **consensus scoring method.** (a-b) We used four consensus scoring methods (all with

956 autodock4, idock, rf\_score, vinardo as input): exp\_z\_score, sum\_rank, best\_rank and

957 worst\_rank. And EF1%, AUC, logAUC, BEDROC ( $\alpha= 80.5$ ) indicators are used to

958 evaluate the screening effect. Panel (a) showed screening results of DUD-E protein

959 structure (red) and AlphaFold protein structure (blue) of 51 targets in DUD-E

960 database. Panel (b) showed screening results of DUD-E protein structure (red),

961 AlphaFold protein structure (blue) and Apo protein structure (green) of 23 targets in

962 DUD-E database. (c) 51 protein targets by exp\_z\_score were used to count the

963 difference of AUC, logAUC, BEDROC ( $\alpha= 321.9$ ,  $\alpha= 80.5$ , and  $\alpha= 20.0$ ) and EF (1%,

964 5%, and 10%) between DUD-E database and AlphaFold protein structure. Negative

965 number (red) indicates that the structure of AlphaFold protein is worse than that of

966 DUD-E protein, and positive number (blue) indicates that the structure of AlphaFold

967 protein is better than that of DUD-E protein. The difference of logAUC sorts the data

968 in the table. There are 20 targets (about 39% of protein targets) whose difference of

969 logAUC is less than -0.03. There are 18 (about 35%) targets whose difference of

970 logAUC is greater than -0.03, and less than or equal to 0.03. There are 13 targets

971 (about 25%) with a difference of logAUC greater than 0.03.

972

973 **Figure 5. Multi-stage screening combined with consensus scoring.** (a) Multi-stage

974 screening flow chart combined with consensus scoring. Prepare small molecules,

975 proteins and search space, and then start multi-stage screening. In the first stage, idock

976 docking is used for virtual screening to obtain idock rank and rf\_score rank, and then

977 ir rank is obtained through the consensus score of idock rank and rf\_score rank. The

978 top x% of the ir rank are taken to enter the second stage. In the second stage, we use

979 the vinardo scoring function to conduct virtual screening to obtain the vinardo rank,

980 and then take the small molecules in the topx% of the ir rank in the idock rank and

981 rf\_score rank and all the small molecules in the vinardo rank for consensus scoring to

982 obtain the irv rank. The top y% of the irv rank are taken to enter the third stage. In the

983 third stage, the autodock4 scoring function is used to conduct virtual screening, and  
984 the autodock4 rank is obtained. Then the small molecules in idock rank, rf\_score rank  
985 and vinardo rank that are y% before irv rank and all small molecules in autodock4  
986 rank are scored by consensus to obtain irva rank. The final ranking is obtained by  
987 ranking the small molecules in irva rank, the small molecules in irv rank that are not  
988 in irva rank, and the small molecules in irv rank that are not in irv rank according to  
989 the original order within the ranking. (b) 51 targets in DUD-E database are in  
990 exp\_z\_score consensus scoring (all with autodock4, idock, rf\_score, vinardo as input).  
991 We integrated multi-stage screening with consensus scoring (plan A, C, E), multi-  
992 stage screening without consensus scoring (plan B, D, F, the flow chart is shown in  
993 Supplementary Information Figure 1) and screening results under a single scoring  
994 function (autodock4, idock, rf\_score, vinardo). The average time of small molecules  
995 of each target in 51 targets (blue, unit: second/cpu/molecule), and the average time of  
996 small molecules of all targets in 51 targets (green, unit: second/cpu/molecule). (c) 51  
997 targets in DUD-E database are in exp\_z\_score Consensus scoring (all with autodock4,  
998 idock, rf\_score, vinardo as input), multi-stage screening with consensus scoring (plan  
999 A, C, E), and multi-stage screening without consensus scoring (plan B, D, F). The  
1000 violin charts demonstrated screening results of EF1%, AUC, logAUC and BEDROC  
1001 ( $\alpha= 80.5$ ). For plan A and plan B, take top 40% in the first stage to enter the second  
1002 stage, and take top 50% in the second stage to enter the third stage. In the first stage of  
1003 plan C and plan D, take top 20% to enter the second stage, and in the second stage,  
1004 take top 50% to enter the third stage. Plan E and plan F take top 10% in the first stage  
1005 to enter the second stage, and take top 50% in the second stage to enter the third stage.

1006

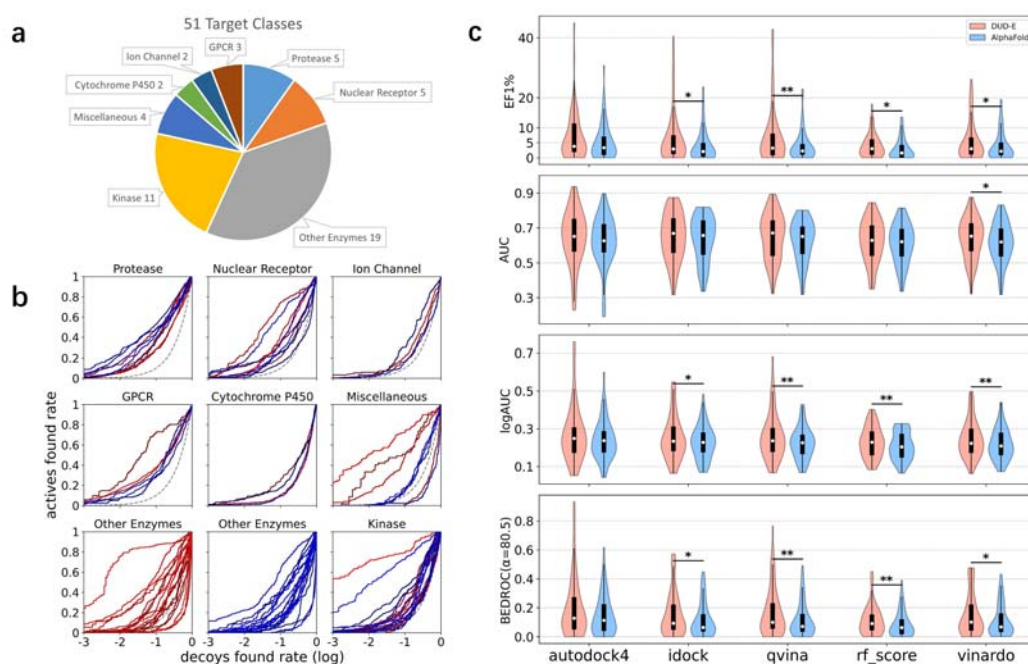
1007 **Figure 6. Relationship between scoring score and hit rate.** (a-j) The hit rate of 51  
1008 targets by Autodock4, idock, rf, vinardo and exp\_z\_score. The red dotted lines  
1009 represent total active small molecules/total small molecules, the rate of each subgraph  
1010 is about 3.0%. The gray dotted lines represent the number of small molecules. The  
1011 green dotted lines represent the total hit rate. The blue dotted lines represent the  
1012 average hit rate of 51 targets. The light gray dotted lines represent the hit rate of each

1013 target of 51 targets. The abscissa of a, b, c, e, f, g, i, j were divided into 22 segments  
 1014 according to the score, of which the end two segments respectively contain data that is  
 1015 greater than or equal to and less than the boundary, and the remaining 20 segments  
 1016 contain data that is left closed and right open within the segment. The abscissa of d  
 1017 and h were divided into 21 segments according to the score. One segment at the right  
 1018 end contains data that is greater than or equal to the boundary, and the other 20  
 1019 segments contain data that is left closed and right open within the segment. A, b, c, d,  
 1020 i are the screening results of protein structure in the DUD-E database, and e, f, g, h, j  
 1021 are AlphaFold protein structure screening results.

1022

## 1023 Figures

1024 Figure 1



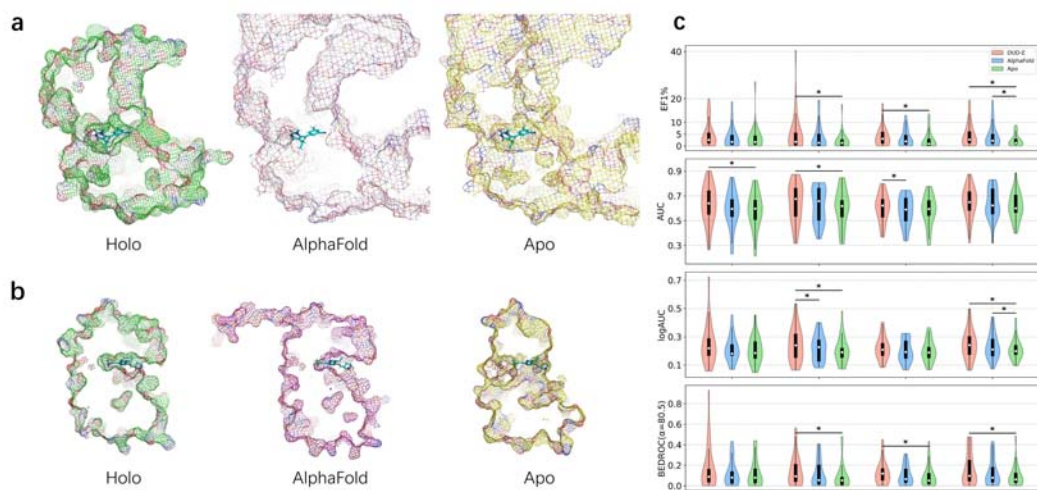
1025

1026

1027

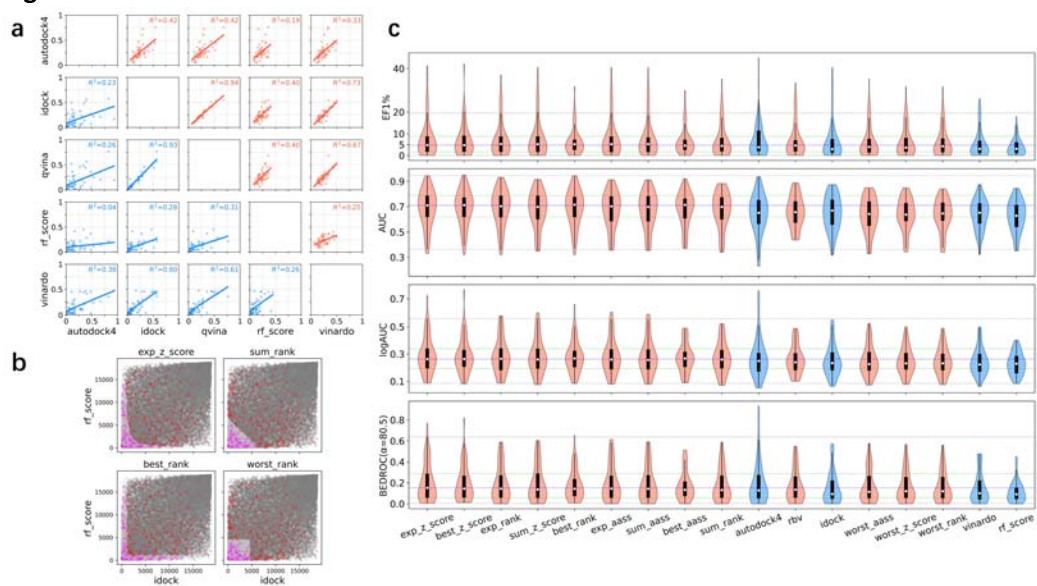
1028

1029 Figure 2



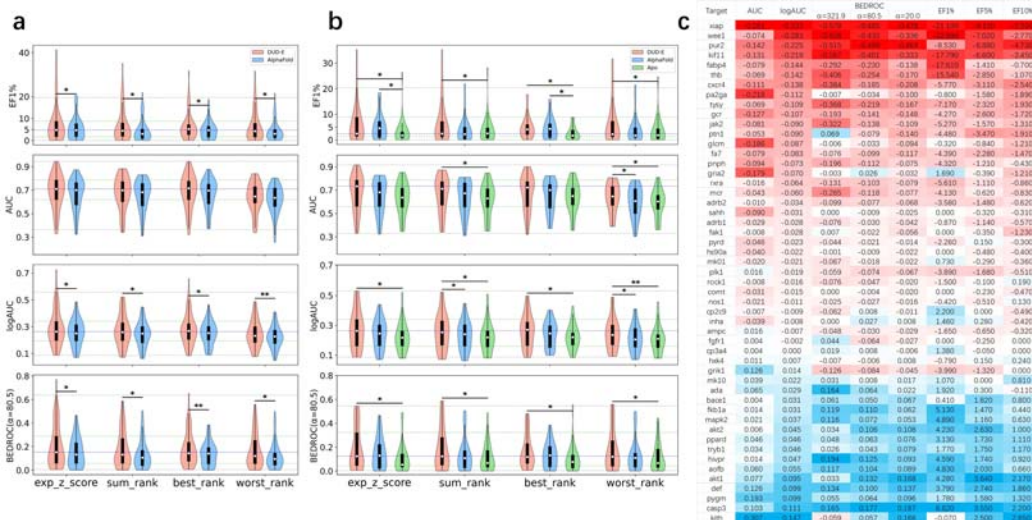
1030  
1031  
1032  
1033  
1034  
1035  
1036

**Figure 3**



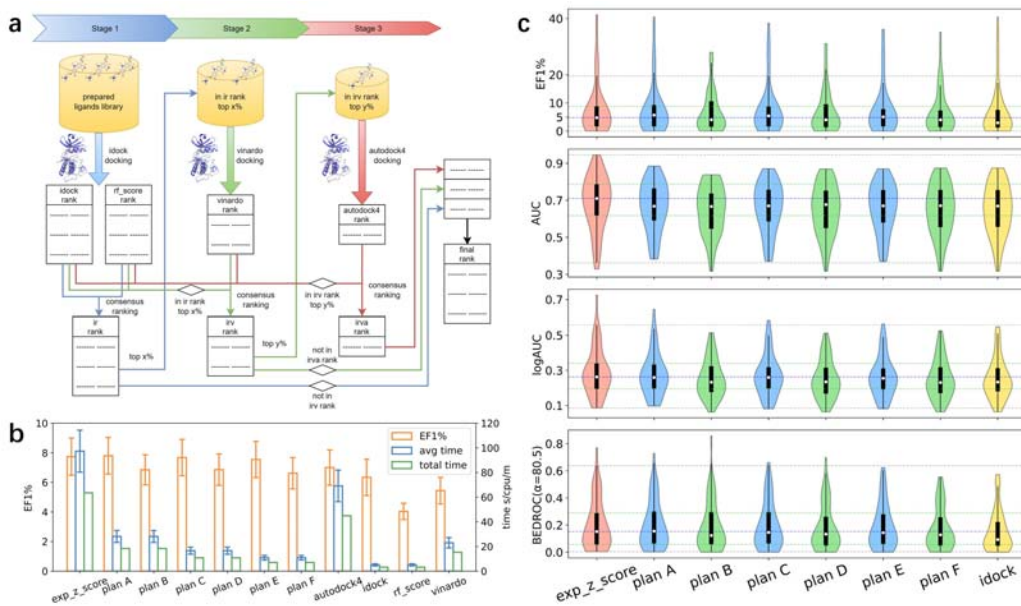
1037  
1038  
1039  
1040  
1041  
1042

**Figure 4**



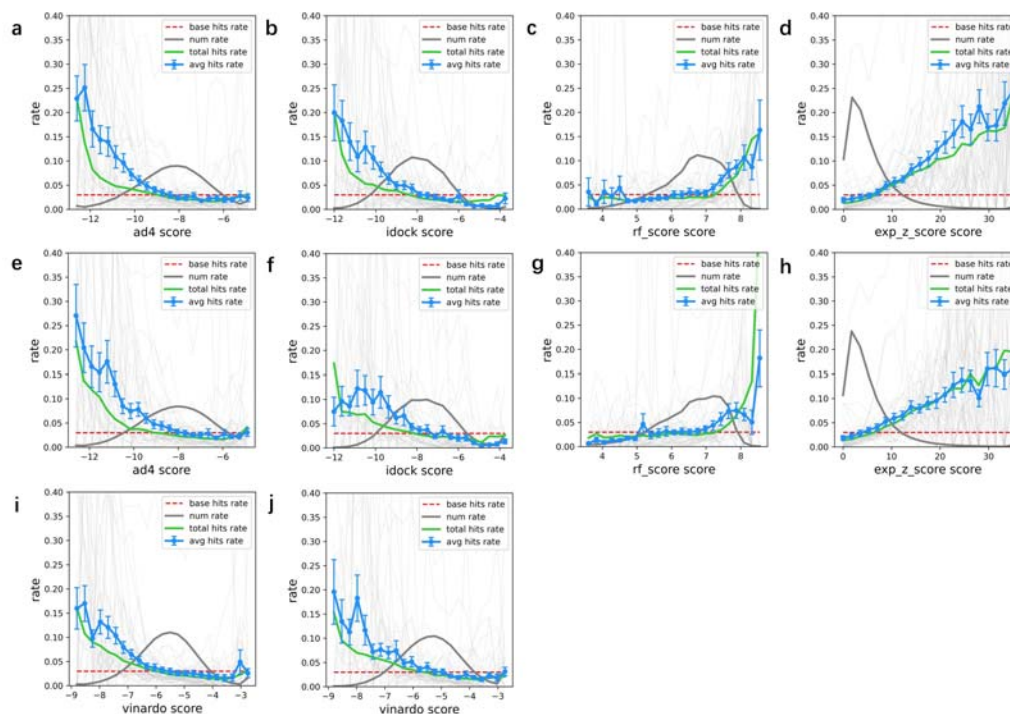
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050

Figure 5



1051  
1052  
1053  
1054  
1055  
1056  
1057

Figure 6



1058

1059

1060

## 1061 Tables

1062 **Table 1. Virtual screening results of protein structure in DUD-E database and**

1063 **AlphaFold.**

	AUC	logAU C	BEDROC( $\alpha$ =321.9)	BEDROC( $\alpha$ =80.5)	BEDROC( $\alpha$ =20.0)	EF1%	EF5%	EF10 %
AF ad4	0.622	0.240	0.183	0.150	0.196	5.284	3.238	2.522
AF idock	0.635	0.228	0.126	0.114	0.166	3.985	2.624	2.226
AF rf_score	0.609	0.208	0.087	0.086	0.140	2.961	2.254	1.984
AF vinardo	0.614	0.217	0.128	0.109	0.157	3.892	2.531	2.105
DUD-E ad4	0.645	0.270	0.229	0.194	0.235	7.002	3.915	2.896
DUD-E idock	0.647	0.252	0.204	0.156	0.199	6.342	3.235	2.574
DUD-E rf_score	0.625	0.228	0.118	0.116	0.172	4.040	2.865	2.336
DUD-E vinardo	0.639	0.245	0.172	0.152	0.198	5.436	3.282	2.509

1064

1065 The average value for 51 targets of AUC, logAUC, BEDROC ( $\alpha= 321.9$ ,  $\alpha= 80.5$ , and

1066  $\alpha= 20.0$ ), EF1%, EF5%, and EF10% were used to evaluate the screening results. AF is

1067 the abbreviation of AlphaFold, indicating the use of AlphaFold protein structure.  
 1068 DUD-E represents the use of the DUD-E database protein structure. Ad4 means using  
 1069 the autodock4 scoring function, idock means using the idock scoring function,  
 1070 rf\_score means using the rf\_score scoring function in the idock software, and vinardo  
 1071 means using the vinardo scoring function.

1072

1073 **Table 2. Virtual screening results of Holo, AlphaFold, Apo protein structure.**

	AUC	logAUC	BEDROC( $\alpha$	BEDROC( $\alpha$	BEDROC( $\alpha$	EF1%	EF5%	EF10%
		C	=321.9)	=80.5)	=20.0)			
AF ad4	0.587	0.212	0.141	0.121	0.167	3.638	2.497	2.100
AF idock	0.625	0.220	0.129	0.115	0.161	3.700	2.291	2.020
AF rf_score	0.576	0.200	0.103	0.101	0.148	3.330	2.439	1.951
AF vinardo	0.633	0.230	0.152	0.126	0.176	4.183	2.622	2.306
Apo ad4	0.580	0.206	0.124	0.122	0.164	3.593	2.396	2.029
Apo idock	0.603	0.202	0.093	0.089	0.139	2.308	2.104	1.797
Apo rf_score	0.581	0.198	0.099	0.092	0.146	2.440	2.249	1.955
Apo vinardo	0.617	0.210	0.074	0.095	0.157	2.311	2.466	2.035
DUD-E ad4	0.628	0.248	0.191	0.172	0.212	4.757	3.267	2.505
DUD-E idock	0.641	0.244	0.179	0.147	0.196	5.455	3.025	2.436
DUD-E rf_score	0.603	0.218	0.134	0.127	0.171	4.016	2.804	2.164
DUD-E vinardo	0.645	0.245	0.171	0.160	0.207	4.667	3.333	2.583

1074

1075 The average value for 51 targets of AUC, logAUC, BEDROC ( $\alpha = 321.9$ ,  $\alpha = 80.5$ , and  
 1076  $\alpha = 20.0$ ), EF1%, EF5%, and EF10% were used to evaluate the screening results. AF is  
 1077 the abbreviation of AlphaFold, indicating the use of AlphaFold protein structure.  
 1078 DUD-E represents the use of the DUD-E database protein structure. Ad4 means using  
 1079 the autodock4 scoring function, idock means using the idock scoring function,  
 1080 rf\_score means using the rf\_score scoring function in the idock software, and vinardo  
 1081 means using the vinardo scoring function.

1082

1083

1084 **Table 3. Evaluation consensus scoring method based on score or ranking.**

	AUC	logAUC	BEDROC( $\alpha$ =321.9)	BEDROC( $\alpha$ =80.5)	BEDROC( $\alpha$ =20.0)	EF1%	EF5%	EF10%
exp_z_score	0.686	0.290	0.248	0.204	0.253	7.736	4.226	3.202
best_z_score	0.687	0.285	0.222	0.193	0.246	7.038	4.162	3.101
exp_rank	0.689	0.284	0.247	0.193	0.244	7.423	4.047	3.104
sum_z_score	0.676	0.282	0.250	0.197	0.243	7.697	4.059	3.068
best_rank	0.688	0.280	0.187	0.180	0.243	6.533	4.138	3.118
exp_aass	0.671	0.279	0.248	0.195	0.240	7.617	3.982	3.021
sum_aass	0.669	0.278	0.248	0.194	0.238	7.506	3.955	2.971
best_aass	0.678	0.272	0.192	0.173	0.230	6.356	3.893	2.998
sum_rank	0.673	0.272	0.234	0.179	0.227	6.782	3.722	2.901
ad4	0.645	0.270	0.229	0.194	0.235	7.002	3.915	2.896
rbv	0.659	0.260	0.211	0.171	0.223	6.735	3.636	3.003
idock	0.647	0.252	0.204	0.156	0.199	6.342	3.235	2.574
worst_aass	0.633	0.250	0.204	0.164	0.206	6.187	3.430	2.608
worst_z_score	0.634	0.249	0.223	0.166	0.204	6.366	3.363	2.534
worst_rank	0.635	0.249	0.223	0.164	0.202	6.258	3.297	2.507
vinardo	0.639	0.245	0.172	0.152	0.198	5.436	3.282	2.509
rf_score	0.625	0.228	0.118	0.116	0.172	4.040	2.865	2.336
p value	0.054	0.075	0.961	0.102	<b>0.028</b>	0.399	0.068	<b>0.028</b>

1085

1086 The average value for 51 targets of AUC, logAUC, BEDROC ( $\alpha= 321.9$ ,  $\alpha= 80.5$ , and  
 1087  $\alpha= 20.0$ ), EF1%, EF5%, and EF10% were used to evaluate the screening results . Ad4  
 1088 means to use the autodock4 scoring function, idock means to use the idock scoring  
 1089 function, rf\_score means to use rf\_score scoring function in idock software, and  
 1090 vinardo means to use vinardo scoring function. In the last line, p-value is the result of  
 1091 paired t-test on exp\_z\_score index data and exp\_rank (ECR) index data. The data  
 1092 ( $p<0.05$ ) has been bold. The remaining data are the screening results of consensus  
 1093 scoring method.

1094



1095

1096 **Table 4. Evaluation of Holo, AlphaFold protein structure using consensus scoring**

1097 **methods.**

	AUC	logAUC	BEDROC( $\alpha$ =321.9)	BEDROC( $\alpha$ =80.5)	BEDROC( $\alpha$ =20.0)	EF1%	EF5%	EF10%
DUD-E	0.686	0.290	0.248	0.204	0.253	7.736	4.226	3.202
AF	0.665	0.256	0.164	0.151	0.208	5.518	3.547	2.744
AF - DUD-E	-0.021	-0.034	-0.084	-0.053	-0.045	-2.218	-0.679	-0.458
p value	0.129	<b>0.014</b>	<b>0.004</b>	<b>0.014</b>	<b>0.024</b>	<b>0.020</b>	0.068	<b>0.043</b>

1098

1099 The screening results of consensus score `exp_z_score` (with `autodock4`, `idock`,  
1100 `rf_score`, `vinardo` as input) show as AUC, logAUC, BEDROC ( $\alpha=321.9$ ,  $\alpha=80.5$ , and  
1101  $\alpha=20.0$ ), EF1%, EF5%, and EF10% on 51 targets. The line named DUD-E means the  
1102 average of 51 structures in DUD-E database. The line named AlphaFold means the  
1103 average of 51 AlphaFold structures. The line named AF – DUD-E means the average  
1104 of 51 AlphaFold structures minus the average of 51 structures in DUD-E database.  
1105 The line named p value is the paired t-test for each index data of DUD-E (Holo) and  
1106 AlphaFold protein structure screening results. The data ( $p<0.05$ ) has been thickened.

1107

1108 **Table 5. Evaluate the protein structure of Holo, AlphaFold and Apo using**  
1109 **consensus scoring methods.**

	AUC	logAUC	BEDROC( $\alpha$ =321.9)	BEDROC( $\alpha$ =80.5)	BEDROC( $\alpha$ =20.0)	EF1%	EF5%	EF10%
DUD-E	0.671	0.271	0.218	0.189	0.237	6.295	3.644	2.918
AF	0.641	0.241	0.169	0.147	0.197	5.161	3.128	2.473
Apo	0.627	0.223	0.103	0.110	0.170	3.223	2.624	2.223
AF - DUD-E	-0.030	-0.030	-0.049	-0.042	-0.040	-1.134	-0.516	-0.445
AF - Apo	0.014	0.018	0.066	0.037	0.026	1.938	0.504	0.249
Apo - DUD-E	-0.044	-0.049	-0.115	-0.079	-0.067	-3.072	-1.020	-0.694
AF DUD-E p value	0.082	0.060	0.192	0.145	0.118	0.294	0.223	0.128
AF Apo p value	0.398	0.243	0.119	0.136	0.233	<b>0.050</b>	0.251	0.364
Apo DUD-E p	0.052	<b>0.020</b>	<b>0.021</b>	<b>0.019</b>	<b>0.029</b>	<b>0.011</b>	0.064	0.055

value											
-------	--	--	--	--	--	--	--	--	--	--	--

1110

1111 The screening results of exp\_z\_score consensus score (with autodock4, idock,  
1112 rf\_score, vinardo as input) show as AUC, logAUC, BEDROC ( $\alpha= 321.9$ ,  $\alpha= 80.5$ , and  
1113  $\alpha= 20.0$ ), EF1%, EF5%, and EF10% on 23 targets. The results of using DUD-E  
1114 database protein structure, AlphaFold (AF) protein structure, Apo protein structure,  
1115 AlphaFold minus DUD-E (AF - DUD-E), AlphaFold minus Apo (AF - Apo), and Apo  
1116 minus DUD-E (Apo – DUD-E) are displayed line by line. P-value is the paired t-test  
1117 for each index data of DUD-E (Holo), AlphaFold, and Apo protein structure screening  
1118 results, respectively. The data of  $p<0.05$  has been thickened.

1119

1120 **Table 6. Early enrichment capacity and time of multi-stage screening combined**  
1121 **with consensus score.**

	exp_z_ score	plan A	plan B	plan C	plan D	plan E	plan F	autod ock4	idock	rf_sco re	vinard o
EF1%	7.74	7.80	6.85	7.68	6.86	7.55	6.62	7.00	6.34	4.04	5.44
avg time	97.34	28.17	28.17	16.65	16.65	10.89	10.89	69.19	5.13	5.13	23.03
total time	63.51	18.38	18.38	10.82	10.82	7.04	7.04	44.93	3.26	3.26	15.31

1122

1123 EF1% indicators of the screening results of 51 targets in the DUD-E database under  
1124 exp\_z\_score consensus score (all with autodock4, idock, rf\_score, vinardo as input),  
1125 multi-stage screening (plan A, C, E) combined with consensus score, multi-stage  
1126 screening (plan B, D, F) without consensus score, and single scoring function  
1127 (autodock4, idock, rf\_score, vinardo). The average small molecule average time of  
1128 each target in 51 targets (avg time, unit: second/cpu/molecule), and the average small  
1129 molecule average time of all targets in 51 targets (total time, unit:  
1130 second/cpu/molecule).

1131

1132