



## 24 **Abstract**

25 A growing number of computational tools have been developed to accurately and rapidly predict  
26 the impact of amino acid mutations on protein-protein relative binding affinities. Such tools have  
27 many applications, for example, designing new drugs and studying evolutionary mechanisms. In  
28 the search for accuracy, many of these methods employ expensive yet rigorous molecular  
29 dynamics simulations. By contrast, non-rigorous methods use less exhaustive statistical  
30 mechanics, allowing for more efficient calculations. However, it is unclear if such methods retain  
31 enough accuracy to replace rigorous methods in binding affinity calculations. This trade-off  
32 between accuracy and computational expense makes it difficult to determine the best method for  
33 a particular system or study. Here, eight non-rigorous computational methods were assessed using  
34 eight antibody-antigen and eight non-antibody-antigen complexes for their ability to accurately  
35 predict relative binding affinities ( $\Delta\Delta G$ ) for 654 single mutations. In addition to assessing  
36 accuracy, we analyzed the CPU cost and performance for each method using a variety of physico-  
37 chemical structural features. This allowed us to posit scenarios in which each method may be best  
38 utilized. Most methods performed worse when applied to antibody-antigen complexes compared  
39 to non-antibody-antigen complexes. Rosetta-based JayZ and EasyE methods classified mutations  
40 as destabilizing ( $\Delta\Delta G < -0.5$  kcal/mol) with high (83-98%) accuracy and a relatively low  
41 computational cost for non-antibody-antigen complexes. Some of the most accurate results for  
42 antibody-antigen systems came from combining molecular dynamics with FoldX with a  
43 correlation coefficient ( $r$ ) of 0.46, but this was also the most computationally expensive method.  
44 Overall, our results suggest these methods can be used to quickly and accurately predict stabilizing  
45 versus destabilizing mutations but are less accurate at predicting actual binding affinities. This  
46 study highlights the need for continued development of reliable, accessible, and reproducible

47 methods for predicting binding affinities in antibody-antigen proteins and provides a recipe for  
48 using current methods.

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

## 70 **Introduction**

71 Protein-protein binding is an essential physiological event that governs a large number of  
72 biological processes in the cell [1]. Amino acid mutations of these proteins can introduce diversity  
73 into genomes, and disrupt or modulate protein-protein interactions by changing the underlying  
74 binding free energy ( $\Delta G$ , i.e. binding affinity), the amount of energy required to form protein  
75 complexes [2]. The binding free energy associated with a protein-protein complex determines the  
76 stability of the complex formation and the conditions for protein-protein association. Accurate  
77 prediction of binding free energies allows us to understand how these affinities can be modified,  
78 and leads to a more comprehensive understanding of protein interactions in living organisms [3].

79

80 Experimental biophysical methods can quantitatively measure change in the protein-protein  
81 binding free energy due to a mutation (i.e. relative binding affinity,  $\Delta\Delta G$ ), but these methods are  
82 typically costly, laborious, and time-consuming since all mutant proteins must be expressed and  
83 purified. Many researchers have developed and utilized computational methods to predict  $\Delta\Delta G$   
84 values for single- or multiple-amino acid mutations (see e.g. [4-6]). Historically, the most  
85 promising in terms of accuracy are rigorous methods based on statistical mechanics that use  
86 molecular dynamics (MD) simulations and thus automatically address conformational flexibility  
87 and entropic effects [7, 8]. However, these methods are computationally expensive since they  
88 employ rigorous sampling and utilize classical mechanics [9] or quantum mechanics [10]  
89 approximations of intermolecular interactions, and require a large number of calculations per time-  
90 step. Because of the expense, rigorous methods are not well-suited to studying large sets of  
91 mutations or large proteins thus necessitating less expensive, non-rigorous methods.

92

93 Non-rigorous high-throughput methods attempt to lower the computational cost, as compared to  
94 rigorous methods, while still providing accurate  $\Delta\Delta G$  predictions. They accomplish this by  
95 including precalculated physico-chemical structural information in combination with predictive  
96 algorithms. The core mechanics that drive these methods fall under numerous classification  
97 umbrellas which have been covered by review articles [11, 12]. These review articles provide a  
98 broad overview but do not provide an unbiased, rigorous, comparative analysis outside of what the  
99 original developers provide. The developers of any given method tend to provide comparisons  
100 with other methods of the same general class to define where their method fits in the current  
101 landscape. BindProfX, for example, is available as a web server and standalone and utilizes  
102 structure-based interface profiles with pseudo counts. Upon release, it was most notably compared  
103 to FoldX (a semi-empirical trained method [13]) and DCOMPLEX (a physics-based method [14])  
104 [15, 16]. iSEE, a statistically trained method based on 31 structure, evolution, and energy-based  
105 terms was tested against FoldX, BindProfX, and BeAtMuSiC (a machine learning-based approach  
106 [17]). Mutabind [18] and some other methods not explored in this work follow a similar testing  
107 methodology [19-21]. While these comparisons are beneficial in providing context for how a given  
108 model fits in the existing research landscape, they are not very robust, since only a narrow subset  
109 of methodologies are included. Conversely for folding stability, Kroncke et al. compared a large  
110 number of available software methods on a small dataset of transmembrane proteins providing a  
111 general overview of performance [6]. Despite the narrow dataset, this study provides a diverse,  
112 useful collection of evaluation metrics between multiple classes of methods. Our intent in this  
113 study is to provide a similar robust comparison of methods for non-rigorous binding affinity  
114 estimation.

115

116 In this work, we evaluate the ability of eight non-rigorous methods to predict relative binding  
117 affinities due to single amino acid mutations. We restrict our study to cases where both an  
118 experimental structure of the complex, and experimentally determined binding affinity values are  
119 available. To investigate the trade-off between speed and accuracy, we chose 16 protein-protein  
120 test complexes with empirical  $\Delta\Delta G$  values for observed mutations. We calculated the  $\Delta\Delta G$  values  
121 for each mutation using all eight methods and compared the results against empirical  $\Delta\Delta G$  values.  
122 The goal of this study was to determine whether software methods that use (most costly) energy  
123 functions with a wider variety of physico-chemical structural features would provide more  
124 accurate binding affinity and interface destabilization predictions compared to those that rely on a  
125 single descriptive (less costly) energy function. We have determined scenarios in which some of  
126 these methods may be better or worse than traditional computational methods in predicting  $\Delta\Delta G$   
127 values.

128

129

130

131

132

133

134

135

136

137

138

## 139 **Methods**

### 140 **Compilation of Experimental $\Delta\Delta G$ Values**

141 To assess the performance of a range of protein-protein binding affinity prediction methods, we  
142 first assembled a dataset containing single amino acid mutations with known experimental  $\Delta\Delta G$   
143 values. This list was assembled from Structural Kinetic and Energetic database of Mutant Protein  
144 Interaction (SKEMPI) version 2.0 [22]. While generating this list, we considered four aspects: (i)  
145 type of protein-protein complex; (ii) availability of quality 3-D structural information; (iii) range  
146 of experimental  $\Delta\Delta G$  values; and (iv) the type of mutations at differing sites on the complex. Our  
147 final dataset contained 654 mutations from 16 protein-protein complexes and their respective  
148 experimental  $\Delta\Delta G$  values. We further categorized these 16 complexes as either non-antibody-  
149 antigen (non-Ab) or antibody-antigen (Ab). Table 1 shows the complexes in our dataset with their  
150 respective non-Ab and Ab categories and the number of mutations associated with each complex.  
151 The dataset contains a total of 401 non-Ab mutations and 253 Ab mutations.

Non-Ab			Ab		
PDB ID	# Mutations	# Residues	PDB ID	# Mutations	# Residues
1a4y [23]	32	583	1bj1[24]	10	547
1brs [25]	30	199	1jrh [26]	42	540
1cbw [27]	31	299	1mlc [28]	11	561
1iar [29]	36	336	1vfb [30]	48	352
1jtg [31]	37	428	1yy9 [32]	16	1058
1lfd [33]	19	254	2jel [34]	43	520
1ppf [35]	190	274	3hfm [36]	71	558
2wpt [37]	26	220	4i77 [38]	12	549

152 **Table 1. Dataset used in our study containing 16 protein complexes.** For both non-Ab (left)  
153 and Ab (right) categories, columns show PDB IDs, total number of residues in a complex, and  
154 number of experimental mutants per complex.

155

## 156 **Selection of Protein-Protein Binding Affinity Methods**

157 Binding affinity prediction methods were chosen to have both a distinct approach to binding  
158 affinity calculation that utilized 3-D structural information and had functional standalone software  
159 in September 2020, available either online or upon request to the author. Table 2 summarizes the  
160 methods selected in this study, their approaches, and their type of scoring functions. For simplicity,  
161 we categorized scoring functions (mathematical functions to calculate  $\Delta\Delta G$  values) as semi-  
162 empirical, statistical, or physics-based. Semi-empirical methods replace as many calculations as  
163 possible with pre-calculated data and are trained using existing crystal structures and known  
164 binding affinity measurements for mutations [39]. Statistical methods use pre-calculated data and  
165 consider changes in coarse structural features such as the change in overall volume [40]. Physics-  
166 based methods use molecular mechanics based-energy functions to estimate enthalpic binding  
167 contributions [14]. In general, statistical or semi-empirical scoring functions involve a training step  
168 where existing datasets are leveraged to determine the weight of input parameters. MD, JayZ, and  
169 EasyE were not developed by training the methods against experimental data designed to improve  
170 predictive power while all other methods utilized this step.

171

172

173

174

175

176

177



Name	Brief Description	Scoring Function	Runtime (CPU hours)
BindProfX [15, 16]	Interface profile score based on conservation of homologous interfaces	Semi-Empirical	1ppf = 0.57 CPUh 1yy9 = 0.73 CPUh
BindProfX(BPX)+FoldX v4 [15, 16]	Profile score weighted and combined with FoldX energy potential	Semi-Empirical	1ppf = 0.62 CPUh 1yy9 = 0.71 CPUh
iSEE [41]	Random forest model using structural, evolutionary, and energy-based features	Statistical	1ppf < 0.01 CPUh 1yy9 < 0.01 CPUh
DCOMPLEX v2 [14]	Structural ideal-gas reference state potential	Physics-Based	1ppf = 0.013 CPUh 1yy9 = 0.001 CPUh
EasyE v1.0 [40, 42]	GMEC-based method utilizing the Rosetta [43, 44] energy function	Statistical	1ppf = 0.48 CPUh 1yy9 = 0.09 CPUh
JayZ v1.0 [40, 42]	Partition-function method utilizing Rosetta energy function	Statistical	1ppf = 0.14 CPUh 1yy9 = 0.21 CPUh
FoldX v4 [13, 39]	Empirical energy score based on various energy parameters (e.g. van der Waals, solvation, electrostatics, hydrogen bonding)	Semi-Empirical	1ppf = 0.42 CPUh 1yy9 = 0.16 CPUh
MD+FoldX v4 [45-47]	Molecular dynamics used to explore conformation space and generate snapshots; FoldX score calculated for each snapshot and averaged	Semi-Empirical	1ppf = 941 CPUh 1yy9 = 4093 CPUh

178 **Table 2. Methods used for comparison in study with a short summary of their approach and**  
179 **scoring function.** Columns (left to right) indicate the method, a brief description of the method,  
180 the type of scoring function used, and runtimes. Runtimes are the amount of CPU hours for  
181 estimating the  $\Delta\Delta G$  for a representative protein complex for Ab (1yy9, 1058 residues) and Non-  
182 Ab (1ppf, 274 residues) categories. Although 1yy9 is roughly four times bigger than 1ppf, the total  
183 runtime may or may not be affected depending on the method used.

184

## 185 Calculation and Comparison of Computational Speed

186 The methods in Table 2 were used to predict  $\Delta\Delta G$  values for each mutation on our experimental  
187 list shown in Table 1. Detailed protocols for predicting  $\Delta\Delta G$  values using each selected method  
188 are provided in the Supplemental Information (see S1 File). Runtimes were determined by using a  
189 representative protein complex from each category: 1ppf, a non-Ab complex with 274 total amino

190 acids, and 1yy9, an Ab complex with 1058 total amino acids (see Table 2). These runtimes are  
191 estimates provided to give a point of comparison between the speeds of different methods. Specific  
192 runtimes will be determined by hardware specifications, method parameters, the number of  
193 mutations being computed, and overall protein size. For MD+FoldX, computational runtime was  
194 the length of time of the MD simulation plus the FoldX runtime for a single mutation. Reporting  
195 runtime in this fashion highlights the large CPUh requirement needed in order to add the sampling  
196 of MD into FoldX calculations. We note that, in contrast to the other methods tested here, the MD  
197 simulations that must be performed for MD+FoldX can be completed very quickly on modern  
198 GPUs, significantly offsetting the high initial cost of the MD+FoldX method. For all other  
199 methods, the algorithms rely either on various pre-calculated data or limited conformational  
200 sampling to calculate  $\Delta\Delta G$  values rapidly.

201

## 202 **Comparing Experimental and Predicted $\Delta\Delta G$ Values**

203 To carry out statistical analysis of our results we built an in-house Python script (see S2 File) that  
204 uses a combination of libraries including matplotlib, numpy, pandas, statistics, scipy, and sklearn.  
205 Using this script, we compared predicted values to experimental  $\Delta\Delta G$  values for each method.

206

207 To evaluate the predictive ability of each method tested, we compared the following correlation  
208 coefficients using our script: concordance ( $\rho_c$ ), Pearson ( $r$ ), Kendall ( $\tau$ ), and Spearman ( $\rho$ ) (see  
209 Table 3). We distinguish between methods that were trained to predict  $\Delta\Delta G$  values from methods  
210 that compute metrics that are expected to linearly correlate with  $\Delta\Delta G$  values. This distinction is  
211 important since for optimal performance we expect a regression line that passes through the  
212 coordinate origin and has a slope of 1, leading to a correlation coefficient equal to 1.

213

Correlation	Brief Description	Type
Concordance	The concordance correlation coefficient ( $\rho_c$ ) measures the degree to which the predicted $\Delta\Delta G$ value equals the actual experimental value (0 indicates no agreement and 1 perfect agreement).	Linear
Pearson	The Pearson correlation coefficient ( $r$ ) measures the degree to which a uniform linear transformation of the predicted $\Delta\Delta G$ values (i.e., a shift and scale change) would yield the actual experimental values (0 indicates no agreement after transformation, 1 perfect agreement, and $-1$ perfect inverse agreement).	Linear
Kendall and Spearman	The rank correlation coefficient measures the degree to which the rank ordering of the predicted $\Delta\Delta G$ values matches the rank ordering of the actual experimental values (0 indicates no agreement after transformation, 1 perfect agreement, and $-1$ perfect inverse agreement). In a normal case, the Kendall correlation ( $\tau$ ) is considered more robust than the Spearman correlation ( $\rho$ ) because of a smaller gross error sensitivity and more efficient due to a smaller asymptotic variance [48].	Rank
AUC and ROC	The receiver operating characteristic (ROC) curve tests several cutoff values for binning mutations as neutral or destabilizing between the most negative calculated $\Delta\Delta G$ value and the most positive calculated $\Delta\Delta G$ value, with true positive rates (sensitivity) calculated at each point. As the true positive rate is calculated, the classifier is moved to less extreme values; this yields the ROC curve. The area under curve (AUC) is a summary statistic that approximates how well the predictor actually discriminates between the two classifications.	N/A

214 **Table 3. Statistical measures used to test the performance of each method in predicting  $\Delta\Delta G$**   
 215 **values.**

216

217 To compare the discriminating power of the methods, we generated receiver operating  
 218 characteristic (ROC) curves (see Table 3). These curves quantify the ability of a method to  
 219 correctly classify point mutations as destabilizing ( $\Delta\Delta G < -0.5$  kcal/mol) or neutral/stabilizing  
 220 ( $\Delta\Delta G > -0.5$  kcal/mol). ROC curves that are skewed toward a higher true positive rate (sensitivity)  
 221 classify mutations more accurately, as quantified by area under curve (AUC, ranging between 1.0  
 222 and 0.5 for perfect and chance classification, respectively).

223

224 We also used our script to parse the results on the basis of several physico-chemical and structural  
225 features to allow us to evaluate the methods based on these characteristics: wild type amino acid  
226 type, mutant amino acid type, protein-protein interacting versus antibody-antigen, secondary  
227 structure classification of the mutation [49, 50], coordination number [51], Sneath index [52],  
228 mostly  $\alpha$ -helical proteins versus mostly  $\beta$ -sheet proteins versus a mix of both  $\alpha$ -helical and  $\beta$ -sheet  
229 proteins, percent exposure, location of the mutation, change in charge, change in polarity, change  
230 in volume, and whether or not the mutation location is predicted as an active or passive residue  
231 [53-55]. The script uses data from S3 File as an input and outputs scatter plots, correlation plots,  
232 receiver operating characteristic (ROC) curves, and box plots to visualize the data, as well as  
233 correlations and standard deviations for each method. All plots in this manuscript were generated  
234 using this script.

235

236

237

238

239

240

241

242

243

244

245

246

## 247 **Results**

248 The purpose of our study was to assess the ability of eight different relative binding affinity  
249 calculation methods (see Table 2) to estimate  $\Delta\Delta G$  values. We selected 16 different protein  
250 complexes (eight Ab, eight non-Ab, see Table 1) with a total of 654 single amino acid mutations.  
251 Each method was then used to estimate  $\Delta\Delta G$  values of 654 mutations and a variety of statistical  
252 measures were employed to assess their predictive ability. We also examined the computational  
253 speed of each method in the context of accuracy to determine its efficiency.

254

### 255 *Non-Antibody-Antigen (non-Ab) Results*

256 Our dataset of eight non-Ab test protein complexes contains 401 total mutations and are mainly  
257 classified as protein-protein systems formed by inhibitors and receptors that range from 199 to 583  
258 residues in size. The distribution and our classification of experimental  $\Delta\Delta G$  values for all non-  
259 Ab test complexes is as follows: 13% of point mutations resulted in  $\Delta\Delta G$  values less than -0.5  
260 kcal/mol (classified as destabilizing); 31% between -0.5 and 0.5 kcal/mol (neutral); and 56%  
261 greater than 0.5 kcal/mol (stabilizing).

262

263 Figures 1 (blue data points and values) and 2 show various performance metrics for each method  
264 to assess their ability to predict the non-Ab  $\Delta\Delta G$  values. Overall, EasyE has the highest correlation  
265 coefficient,  $r = 0.62$ , and iSEE has the lowest,  $r = 0.17$  (see Figures 1 and 2). JayZ and EasyE,  
266 both of which utilize Rosetta's conformational sampling algorithms, consistently have the best  $r$   
267 values for non-Ab mutations.

268

269

270 **Figure 1. Calculated  $\Delta\Delta G$  values (x-axis) compared to experimental  $\Delta\Delta G$  values (y-axis) for**  
271 **each method tested in this study.** Black, red, and blue lines are simple linear regressions from  
272 which  $r$  are derived. The red points are a scatter for Ab complexes and the blue points are for non-  
273 Ab complexes. The dashed line is the  $y = x$  line measuring perfect agreement between predicted  
274 and experimental  $\Delta\Delta G$  values. The solid black, red, and blue lines indicate a linear relationship  
275 between calculated and experimental observations for all data points, Ab complexes, and non-Ab  
276 complexes respectively. The top values in black, red, and blue match the root-mean-square error  
277 and the bottom values indicate  $r$  for all values, Ab values, and non-Ab values respectively.  
278

279 **Figure 2. Performance of each method for non-Ab complexes (401 total mutations) in**  
280 **predicting true  $\Delta\Delta G$  values ( $\rho_c$ ), linearly correlated  $\Delta\Delta G$  values ( $r$ ), and rank order ( $\rho$  and**  
281  **$\tau$ ).** The error for each method is reported under the correlation points.  
282

283  
284 Figure 3 shows the ROC plot for all the tested methods. These ROC plots highlight how well a  
285 method can discriminate between stabilizing and destabilizing mutations. JayZ (0.84), EasyE  
286 (0.83), DCOMPLEX (0.82), FoldX (0.79), and MD+FoldX (0.76) have the highest AUC.  
287 Combined with the results from Figures 1 and 2, for the systems studied here, JayZ and EasyE  
288 methods are the best overall performers in terms of accuracy, discriminating stabilizing mutations  
289 from destabilizing, and ranking mutations based on their experimental  $\Delta\Delta G$  values.

290  
291 **Figure 3. Receiver operating characteristic (ROC) curves for non-Ab complexes of the**  
292 **classification of variants as stabilizing ( $\Delta\Delta G < -0.5$  kcal/mol) or destabilizing ( $\Delta\Delta G > 0.5$**   
293 **kcal/mol).** The values in the legend represent the area-under-curve (AUC). The higher the value,  
294 the better method is at discriminating between destabilizing and destabilizing mutations.  
295

296  
297 Table 2 reports CPUh required (i.e. runtimes) for each method to calculate  $\Delta\Delta G$  for the entire list  
298 of mutations for a representative non-Ab protein complex. BindProfX, BindProfX(BPX)+FoldX,  
299 JayZ, and EasyE allow users to specify a list of mutations that the method is then able to calculate  
300 in one setting. This list can be optimized based on the available hardware to achieve efficiency.

301 iSEE requires significant preparatory work (see File S1) prior to calculation, but once completed,  
302 it calculates the  $\Delta\Delta G$  values for the entire list of mutations nearly instantly. DCOMPLEX is not  
303 as flexible out of the box but can handle large numbers of mutations through an automated script.  
304 For MD+FoldX, 1yy9 (roughly four times larger than 1ppf) requires considerably more CPUh to  
305 calculate. All other methods calculate 1yy9 in a shorter time frame than 1ppf. This may seem  
306 counterintuitive. However, MD must statistically sample the conformational energy of the entire  
307 complex, while all other methods use algorithms that are likely impacted more by the number of  
308 residues involved in the interaction rather than the protein size. Overall, DCOMPLEX has a much  
309 faster runtime compared to other methods, and if the goal is to determine stabilizing and  
310 destabilizing non-Ab mutations, it offers similar discriminating power to JayZ and EasyE, at a  
311 fraction of the computational cost. JayZ estimates  $\Delta\Delta G$  value of one mutation in  $\sim 2.7$  s, EasyE in  
312  $\sim 9.1$  s, but DCOMPLEX requires just  $\sim 0.25$  s. Overall, EasyE appears to be the best option for  
313 balancing accuracy and speed and DCOMPLEX is recommended for discriminating between  
314 stability and destabilizing mutations.

315  
316 A method might not be a good overall performer in predicting  $\Delta\Delta G$  values but could still perform  
317 well for mutations with certain physico-chemical and structural features. Therefore, we calculated  
318 various statistical measures to assess each method on unique subsets of mutations (see Table 4 and  
319 SI Figs S1-4). This table shows eight different data subsets with two  $r$  per method. EasyE has the  
320 highest  $r$  for non-Ab for five out of eight subsets (wild type non-gly or non-pro, alpha helix, beta  
321 sheet, surface exposure, and large volume changes). Where this method did not have the highest  
322  $r$ , it had either the second or third highest  $r$ . JayZ mirrors the performance of EasyE in all the same  
323 categories and performs better than Easy in the neutral charge subset. These results further

324 highlight the versatility of EasyE’s and JayZ’s performance in estimating the effects of non-Ab  
 325 mutations compared to the other methods tested in this study. All methods apart from iSEE and  
 326 BindProfX perform surprisingly well in the WT Gly or Pro subset. iSEE’s performance in this  
 327 subset, while still mediocre compared to the other tested methods, is substantially better than in all  
 328 other subsets.  
 329

Method	WT Gly or Pro	WT Non-Gly or Non-Pro	Alpha Helix	Beta Sheet	Surface Exposure	Neutral Charge	Hydrophobic to Polar	Large Vol Changes
BindProfX	Non-Ab: 0.11 Ab: -0.03	Non-Ab: 0.33 Ab: 0.23	Non-Ab: 0.29 Ab: 0.16	Non-Ab: 0.29 Ab: <b>0.52</b>	Non-Ab: 0.22 Ab: 0.09	Non-Ab: <b>0.37</b> Ab: 0.28	Non-Ab: 0.33 Ab: 0.17	Non-Ab: 0.13 Ab: 0.42
BPX+FoldX	Non-Ab: <b>0.81</b> Ab: 0.09	Non-Ab: 0.45 Ab: 0.34	Non-Ab: 0.43 Ab: 0.39	Non-Ab: 0.43 Ab: <b>0.54</b>	Non-Ab: 0.32 Ab: 0.21	Non-Ab: 0.52 Ab: 0.41	Non-Ab: 0.41 Ab: 0.26	Non-Ab: 0.71 Ab: <b>0.50</b>
FoldX	Non-Ab: <b>0.85</b> Ab: -0.11	Non-Ab: 0.45 Ab: 0.25	Non-Ab: 0.39 Ab: 0.25	Non-Ab: 0.39 Ab: 0.31	Non-Ab: 0.50 Ab: 0.26	Non-Ab: 0.42 Ab: <b>0.41</b>	Non-Ab: 0.41 Ab: 0.11	Non-Ab: 0.63 Ab: -0.32
MD+FoldX	Non-Ab: <b>0.83</b> Ab: <b>0.71</b>	Non-Ab: 0.49 Ab: <b>0.42</b>	Non-Ab: 0.44 Ab: <b>0.54</b>	Non-Ab: 0.44 Ab: 0.49	Non-Ab: 0.47 Ab: 0.35	Non-Ab: 0.46 Ab: 0.46	Non-Ab: <b>0.46</b> Ab: <b>0.31</b>	Non-Ab: 0.71 Ab: 0.35
DCOMPLEX	Non-Ab: <b>0.65</b> Ab: <b>0.89</b>	Non-Ab: 0.34 Ab: 0.37	Non-Ab: 0.33 Ab: 0.31	Non-Ab: 0.33 Ab: 0.30	Non-Ab: 0.52 Ab: 0.27	Non-Ab: 0.36 Ab: <b>0.56</b>	Non-Ab: 0.38 Ab: 0.16	Non-Ab: 0.62 Ab: 0.28
JayZ	Non-Ab: 0.80 Ab: <b>0.54</b>	Non-Ab: 0.49 Ab: 0.24	Non-Ab: 0.44 Ab: -0.06	Non-Ab: 0.44 Ab: 0.16	Non-Ab: 0.59 Ab: <b>0.36</b>	Non-Ab: <b>0.62</b> Ab: 0.26	Non-Ab: 0.41 Ab: 0.01	Non-Ab: <b>0.83</b> Ab: 0.19
EasyE	Non-Ab: 0.80 Ab: 0.29	Non-Ab: <b>0.51</b> Ab: 0.22	Non-Ab: <b>0.51</b> Ab: 0.06	Non-Ab: <b>0.51</b> Ab: 0.03	Non-Ab: <b>0.60</b> Ab: <b>0.35</b>	Non-Ab: 0.61 Ab: 0.23	Non-Ab: 0.45 Ab: 0.02	Non-Ab: <b>0.84</b> Ab: 0.18
iSEE	Non-Ab: <b>0.43</b> Ab: -0.43	Non-Ab: 0.28 Ab: -0.16	Non-Ab: 0.05 Ab: -0.04	Non-Ab: 0.05 Ab: -0.24	Non-Ab: 0.15 Ab: <b>0.11</b>	Non-Ab: 0.15 Ab: -0.11	Non-Ab: 0.14 Ab: -0.02	Non-Ab: 0.24 Ab: -0.44

330 **Table 4. All methods  $r$  with respect to certain subsets.** “WT Gly or Pro” are wild type amino  
 331 acids that are either glycine or proline. “WT Non-Gly or Non-Pro” are wild type amino acids that  
 332 are neither glycine nor proline. “Alpha Helix” are mutations that occur in a helix structure. “Beta  
 333 Sheet” are mutations that occur in a beta structure. “Surface Exposure” are mutations that occur in  
 334 an amino acid that have relative solvent accessibility values between 0 and 10%. “Neutral Charge”  
 335 is a neutrally charged wild type amino acid mutating to a neutrally charged mutant amino acid.  
 336 “Hydrophobic to Polar” is a hydrophobic or polar wild type amino acid mutating to a polar or  
 337 hydrophobic mutant amino acid, respectively. “Larger Vol Changes” is a mutant amino acid that  
 338 is greater than 40% larger than the wild type amino acid. Values that are bolded are the highest  $r$   
 339 for each method and protein type. Values that are red or blue are the highest  $r$  for each subset, blue  
 340 for non-Ab and red for Ab.  
 341

#### 342 *Antibody-Antigen (Ab) Results*

343 Our dataset of eight Ab test protein complexes contains 253 mutations and the proteins range in  
 344 size from 352 to 1058 residues. The distribution and our classification of experimental  $\Delta\Delta G$  values  
 345 for all Ab test complexes is as follows: 5% of point mutations resulted in  $\Delta\Delta G$  values less than -



346 0.5 kcal/mol (classified as destabilizing); 40% between -0.5 and 0.5 kcal/mol (neutral); and 55%  
347 greater than 0.5 kcal/mol (stabilizing).

348  
349 Figures 1 (data points and values in red), 4, and 5 show the performance of each method in  
350 predicting the  $\Delta\Delta G$  values of Ab mutations. Overall, the highest correlation is for MD+FoldX with  
351  $r = 0.39$  and the lowest is iSEE with  $r = -0.09$  (see Figures 1 and 4). An interesting trend is that  
352 the methods with the highest  $r$  values for non-Ab complexes do not have the highest  $r$  for Ab  
353 complexes.

354  
355 **Figure 4. Performance of each evaluated method for Ab complexes (253 total mutations) in**  
356 **predicting true  $\Delta\Delta G$  values ( $\rho_c$ ), linearly correlated  $\Delta\Delta G$  values ( $r$ ), and rank order ( $\rho$  and**  
357  **$\tau$ ).** The error for each method is reported under the correlation points.  
358

359 **Figure 5. Receiver operating characteristic curves of the classification of variants that are**  
360 **more destabilized or less destabilized than 0.5 kcal/mol.** The values in the legend represent the  
361 area-under-curve (AUC). The higher the value, the better the prediction capability of the method.  
362

363  
364 Figure 5 shows the ROC plot for all the tested Ab methods. These ROC plots highlight how well  
365 a method is actually able to discriminate between stabilizing and destabilizing mutations.  
366 Compared to non-Ab complexes, all methods performed better for antibody-antigen complexes  
367 except for FoldX and DCOMPLEX which were marginally worse. JayZ (0.97), EasyE (0.98),  
368 FoldX (0.85), and MD+FoldX (0.82) had the highest AUC values. Combined with the results from  
369 Figures 1 and 4, at least for the systems studied here, it appears that the MD+FoldX method is the  
370 best overall performer in terms of accuracy, discriminating stabilizing mutations from  
371 destabilizing, and ranking mutations based on their experimental  $\Delta\Delta G$  values.

372

373 Compared to other methods, EasyE has a much faster runtime and is recommended if the goal is  
374 to discriminate between stabilizing and destabilizing ( $\Delta\Delta G$  for one mutation takes  $\sim 21$  s, see Table  
375 2). By comparison, MD+FoldX cost  $\sim 941$  CPUh for one mutation of 1yy9. DCOMPLEX provides  
376 a slightly lower  $r$  (0.31) and computational cost ( $\sim 0.35$  s) for one mutation of 1yy9. Overall,  
377 MD+FoldX appears to be the best option for accuracy and EasyE or JayZ are the best options for  
378 discriminating between destabilizing and stabilizing mutations.

379

380 Table 4 summarizes the ability of each method to predict  $\Delta\Delta G$  values for subsets of Ab mutations.  
381 Most methods had mediocre  $r$  values less than 0.60. The exceptions to this are MD+FoldX and  
382 DCOMPLEX within the WT Gly or Pro subset with  $r = 0.71$  and  $0.89$ , respectively. MD+FoldX  
383 has the highest  $r$  for Ab complexes for three of the eight subsets (WT nonGly or nonPro, alpha  
384 helix, and hydrophobic to polar). BPX+FoldX has the highest  $r$  in two of the eight subsets (beta  
385 sheet and large volume changes). For the beta sheet subset, BindProfX had the second highest  $r$ .  
386 DCOMPLEX had the highest  $r$  for two different subsets (WT Gly or Pro and neutral charge). In  
387 the surface exposure subset, JayZ and EasyE both had nearly identical  $r$  ( $0.36$  and  $0.35$   
388 respectively), the highest for this subset, but substantially worse than they did for non-Ab  
389 complexes.

390

391

392

393

394

## 395 Discussion

396 We assessed the performance of eight distinct protein-protein binding affinity calculation methods  
397 that use 3-D structural information. To test the performance of these methods, we selected 16  
398 different protein complexes (see Table 1) with a total of 654 single amino acid mutations: eight  
399 antigen-antibody complexes (Ab, 253 mutations) and eight non-antigen-antibody (Non-Ab, 401  
400 mutations) complexes. Each method was used to estimate  $\Delta\Delta G$  values of the 654 mutations, a  
401 variety of statistical measures, CPU cost, and physico-chemical structural features to assess the  
402 performance. Our results suggest each method has both strengths and weaknesses depending on  
403 the properties of the protein system. Most methods did not perform well when applied to mutations  
404 in Ab complexes compared to non-Ab complexes. Rosetta-based JayZ and EasyE were able to  
405 classify mutations as destabilizing ( $\Delta\Delta G < -0.5$  kcal/mol) with high (83-98%) accuracy at  
406 relatively low computational cost. Some of the best results for Ab systems came from combining  
407 MD simulations with FoldX with a  $r$  coefficient of 0.39, but at the highest computational cost of  
408 all the tested methods.

409  
410 Figure 1 summarizes the performance of each method in terms of its ability to estimate  $\Delta\Delta G$  values  
411 for all (non-Ab + Ab) single mutations. None of the test methods show a very high  $r$  between  
412 experimental and predicted  $\Delta\Delta G$  values. Two of the best performing methods, JayZ and EasyE,  
413 both have an  $r$  of 0.49 for all mutations, with a higher  $r$  of 0.61 and 0.62 respectively for non-Ab  
414 complexes. These results agree with published results from the authors of JayZ and EasyE. Our  
415 results agree moderately with published results from iSEE (they obtained  $r = 0.25$ , we obtained  $r$   
416  $= 0.17$ ) and BindProfX (they used a much larger dataset). Published results for DCOMPLEX show  
417 a very good correlation of  $r = 0.87$ ; much larger than what we obtained here. This difference is

418 very likely due to the dataset size and compilation; DCOMPLEX was originally tested against 69  
419 experimental data points, compared to the 654 values used here. MD+FoldX has an  $r$  of 0.39 for  
420 Ab complexes and appears to perform well for larger systems, which could indicate the importance  
421 of conformational sampling for antibody-antigen systems. Other methods used in this study have  
422 little to no conformational sampling which could explain their poor performance on Ab complexes.  
423 By contrast, these same methods perform well for non-Ab complexes, suggesting that  
424 conformational sampling is not the limiting factor to achieve accurate results for these protein  
425 complexes. For example, FoldX has a trained scoring function derived using a dataset of mostly  
426 non-Ab complexes and performs poorly for Ab complexes when using a single structure (see Table  
427 2). However, when used with snapshots from an MD simulation, this same method outperforms  
428 all other methods selected in this study. This highlights the need for conformational sampling for  
429 reliable and efficient predictions of binding affinity for some systems. In our previous study, we  
430 combined coarse-grained forcefield with umbrella sampling to calculate  $\Delta\Delta G$  values for eight  
431 mutations of 3hfm Ab complex (one of the test systems in this study) and obtained better  
432 predictions than FoldX and MD+FoldX [56]. This study further emphasizes the need for better  
433 conformational strategies for some systems.

434  
435 Statistical measures used to analyze performance are listed and defined in Table 3. For Ab,  
436 BPX+FoldX, MD+FoldX, and DCOMPLEX have the highest  $r$  values of the methods in our study  
437 (see Figure 4). MD+FoldX appears to be the most accurate method for Ab complexes. BindProfX,  
438 FoldX, JayZ, EasyE, and iSEE have low  $r$  and  $\rho_c$  indicating that affinities estimated using these  
439 methods do not correlate well with experimental  $\Delta\Delta G$  values using a linear transformation. Also,

440 the  $\tau$  and  $\rho$  were lower compared to MD+FoldX, indicating these methods do poorly at calculating  
441 a rank order that matches experimental data.

442  
443 The ROC curves allow us to determine each method's ability to classify mutations as either  
444 destabilizing or neutral/stabilizing (Figures 3 and 5). For non-Ab complexes, JayZ (0.84 AUC)  
445 and EasyE (0.83 AUC) have the best true positive rate followed by DCOMPLEX (0.82 AUC). For  
446 Ab complexes, JayZ (0.97 AUC) and EasyE (0.98 AUC) have better true positive rates than  
447 MD+FoldX, the method with the highest  $r$  value. If classification of destabilizing vs stabilizing is  
448 the primary need, then JayZ or EasyE are both recommended over the other methods tested here  
449 due to their high accuracy and fast runtime.

450  
451 While accuracy is generally the main reason for choosing a particular method, computational  
452 efficiency is also an important consideration, especially when predicting the effects of a large  
453 number of mutations. Here, we discuss the performance of each method in terms of its trade-off  
454 between speed and accuracy for predicting  $\Delta\Delta G$  values. For all single mutations and our non-Ab  
455 subset, EasyE and JayZ performed well; JayZ is the faster method of the two with EasyE at a  
456 similar speed to FoldX. DCOMPLEX is more accurate than FoldX for all single mutations and has  
457 similar accuracy as FoldX for non-Ab mutations, but at much lower cost. MD+FoldX has similar  
458 accuracy to DCOMPLEX for all single mutations and has similar accuracy to FoldX in non-Ab  
459 mutations but is by far the most computationally expensive method we tested. Although a  
460 synergistic combination of BPX+FoldX implements several structural and physico-chemical  
461 interaction terms in its algorithm, computation time was longer than all but MD+FoldX without a  
462 concomitant improvement in  $r$ . We note that this method is perhaps the most accessible of those

463 tested, due to the easy-to-use online server interface and accuracy that is similar to FoldX for most  
464 systems. BindProfX utilizes the same scoring profile as BPX+FoldX without the FoldX  
465 calculations. In this case, accuracy decreased while calculation speed remained similar to  
466 BPX+FoldX. iSEE, the least correlating method, employs the widest variety of information to  
467 obtain relative binding affinity predictions and is the fastest of all methods (not including the non-  
468 trivial preparation time). For Ab complexes, MD+FoldX, the slowest of all the methods, had the  
469 highest accuracy, followed by DCOMPLEX. iSEE is again the fastest of all methods but also the  
470 least accurate. BindProfX utilizes several pre-calculated physico-chemical structural data in its  
471 scoring function while, JayZ and EasyE each layer an additional predictive calculating feature on  
472 top of Rosetta's backbone sampling, adding complexity to the predictive algorithms. However, all  
473 three have similar  $r$  yet they do not achieve the accuracy of MD+FoldX. Overall, for non-Ab  
474 complexes, EasyE and JayZ appear to have the best balance between speed and accuracy of the  
475 methods we tested. For Ab complexes, DCOMPLEX appears to have the best balance.

476  
477 We have demonstrated that all the tested methods have specific strengths and weaknesses; some  
478 perform better in specific contexts (Table 4), and some have longer runtimes to obtain similar  
479 predictive power to comparably faster methods. This highlights the complexity of the physico-  
480 chemical properties and structural features that drive, and limit, these predictive models. Our  
481 results can be used to make informed decisions for methods that may be preferable for a particular  
482 study or system. Table 4 suggests that if the goal is to estimate only the order of magnitude or sign  
483 of relative binding affinities, then the preferred method will likely be very different than if the goal  
484 is to obtain the best possible accuracy for antibody-antigen systems. To improve accessibility, we  
485 have generated an in-house Python script (provided in the supplement with the full dataset used in

486 this work) that can be used to parse any of the parameters used in this study and provide tailored  
487 information. This information in combination with the runtime and other details provided in this  
488 study can be used to inform users on methods that can provide the best accuracy and efficiency for  
489 a given protein-protein complex type, set of physico-chemical features or structural parameters,  
490 and set of mutations. Additionally, the script can be extended to other methods and feature-sets,  
491 potentially elucidating specific problems or areas of improvement to existing and future methods.

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

## 507 **Conclusions**

508 In this study, we have assessed the accuracy and efficiency of eight computational methods on  
509 predicting binding affinity changes due to single amino acid mutations. Methods were tested on  
510 16 different protein complexes: eight antigen-antibody (Ab) and eight non-antigen-antibody (Non-  
511 Ab) complexes. While some methods perform consistently better than others, how well each  
512 performs depends on the physico-chemical and structural components of each complex. EasyE  
513 was the most accurate for non-Ab complexes, and MD+FoldX was most accurate for Ab  
514 complexes. JayZ and EasyE were better able to distinguish between destabilizing ( $\Delta\Delta G > 0.5$   
515 kcal/mol) and stabilizing ( $\Delta\Delta G < -0.5$  kcal/mol) as compared to any other method. Future work  
516 could include more systems or different methods, including those that are solely web server-based  
517 in order to expand and better refine our conclusions on their predictive capability.

518

519

520

521

522

523

524

525

526

527

528

529



## 530 References

- 531
- 532 1. Jones S, Thornton JM. Principles of protein-protein interactions. *Proceedings of the National*  
533 *Academy of Sciences*. 1996;93(1):13-20.
  - 534 2. Yates CM, Sternberg MJE. The Effects of Non-Synonymous Single Nucleotide Polymorphisms  
535 (nsSNPs) on Protein-Protein Interactions. *Journal of Molecular Biology*. 2013;425(21):3949-63.
  - 536 3. Baaden M, Marrink SJ. Coarse-grain modelling of protein-protein interactions. *Curr Opin Struct*  
537 *Biol*. 2013;23(6):878-86. PubMed PMID: 24172141.
  - 538 4. Ezkurdia I, Bartoli L, Fariselli P, Casadio R, Valencia A, Tress ML. Progress and challenges in  
539 predicting protein-protein interaction sites. *Briefings in bioinformatics*. 2009;10(3):233-46. PubMed  
540 PMID: 19346321.
  - 541 5. Kastritis PL, Bonvin AM. Are scoring functions in protein-protein docking ready to predict  
542 interactomes? Clues from a novel binding affinity benchmark. *J Proteome Res*. 2010;9(5):2216-25.  
543 PubMed PMID: 20329755.
  - 544 6. Kroncke BM, Duran AM, Mendenhall JL, Meiler J, Blume JD, Sanders CR. Documentation of an  
545 Imperative To Improve Methods for Predicting Membrane Protein Stability. *Biochemistry*.  
546 2016;55(36):5002-9.
  - 547 7. Bernardi RC, Melo MC, Schulten K. Enhanced sampling techniques in molecular dynamics  
548 simulations of biological systems. *Biochim Biophys Acta*. 2015;1850(5):872-7. PubMed PMID: 25450171;  
549 PubMed Central PMCID: PMC4339458.
  - 550 8. Spiwok V, Sucur Z, Hosek P. Enhanced sampling techniques in biomolecular simulations.  
551 *Biotechnol Adv*. 2015;33(6 Pt 2):1130-40. PubMed PMID: 25482668.
  - 552 9. Gumbart JC, Roux B, Chipot C. Efficient Determination of Protein-Protein Standard Binding Free  
553 Energies from First Principles. *J Chem Theory Comput*. 2013;9(8):3789-98.
  - 554 10. Pokorná P, Kruse H, Krepl M, Šponer J. QM/MM Calculations on Protein-RNA Complexes:  
555 Understanding Limitations of Classical MD Simulations and Search for Reliable Cost-Effective QM  
556 Methods. *J Chem Theory Comput*. 2018;14(10):5419-33.
  - 557 11. Geng C, Xue LC, Roel-Touris J, Bonvin AMJJ. Finding the  $\Delta\Delta G$  spot: Are predictors of binding  
558 affinity changes upon mutations in protein-protein interactions ready for it? *WIREs Computational*  
559 *Molecular Science*. 2019;9(5):e1410.
  - 560 12. Gromiha MM, Yugandhar K, Jemimah S. Protein-protein interactions: scoring schemes and  
561 binding affinity. *Curr Opin Struct Biol*. 2016;44:31-8. PubMed PMID: 27866112.
  - 562 13. Schymkowitz J, Borg J, Stricher F, Nys R, Rousseau F, Serrano L. The FoldX web server: an online  
563 force field. *Nucleic acids research*. 2005;33(Web Server issue):W382-8. PubMed PMID: 15980494;  
564 PubMed Central PMCID: PMC1160148.
  - 565 14. Liu S, Zhang C, Zhou H, Zhou Y. A physical reference state unifies the structure-derived potential  
566 of mean force for protein folding and binding. *Proteins: Structure, Function, and Bioinformatics*.  
567 2004;56(1):93-101.
  - 568 15. Xiong P, Zhang C, Zheng W, Zhang Y. BindProfX: Assessing Mutation-Induced Binding Affinity  
569 Change by Protein Interface Profiles with Pseudo-Counts. *J Mol Biol*. 2017;429(3):426-34.
  - 570 16. Brender JR, Zhang Y. Predicting the Effect of Mutations on Protein-Protein Binding Interactions  
571 through Structure-Based Interface Profiles. *PLoS Comp Biol*. 2015;11(10):e1004494. PubMed PMID:  
572 26506533; PubMed Central PMCID: PMC4624718.
  - 573 17. Dehouck Y, Kwasigroch JM, Rooman M, Gilis D. BeAtMuSiC: Prediction of changes in protein-  
574 protein binding affinity on mutations. *Nucleic Acids Res*. 2013;41(Web Server issue):W333-9. PubMed  
575 PMID: 23723246; PubMed Central PMCID: PMC3692068.
  - 576 18. Li M, Simonetti FL, Goncarenco A, Panchenko AR. MutaBind estimates and interprets the effects  
577 of sequence variants on protein-protein interactions. *Nucleic Acids Res*. 2016;44(W1):W494-W501.
  - 578 19. Vreven T, Hwang H, Pierce BG, Weng Z. Prediction of protein-protein binding free energies.  
579 *Protein Sci*. 2012;21(3):396-404.

- 580 20. Rodrigues CHM, Myung Y, Pires DEV, Ascher DB. mCSM-PPI2: predicting the effects of  
581 mutations on protein–protein interactions. *Nucleic Acids Res.* 2019;47(W1):W338-W44.
- 582 21. Jemimah S, Sekijima M, Gromiha MM. ProAffiMuSeq: sequence-based method to predict the  
583 binding free energy change of protein–protein complexes upon mutation using functional classification.  
584 *Bioinformatics.* 2019;36(6):1725-30.
- 585 22. Jankauskaitė J, Jiménez-García B, Dapkūnas J, Fernández-Recio J, Moal IH. SKEMPI 2.0: an  
586 updated benchmark of changes in protein–protein binding energy, kinetics and thermodynamics upon  
587 mutation. *Bioinformatics.* 2018;35(3):462-9.
- 588 23. Papageorgiou AC, Shapiro R, Acharya KR. Molecular recognition of human angiogenin by  
589 placental ribonuclease inhibitor—an X-ray crystallographic study at 2.0 Å resolution. *The EMBO Journal.*  
590 1997;16(17):5162-77.
- 591 24. Muller YA, Chen Y, Christinger HW, Li B, Cunningham BC, Lowman HB, et al. VEGF and the  
592 Fab fragment of a humanized neutralizing antibody: crystal structure of the complex at 2.4 Å  
593 resolution and mutational analysis of the interface. *Structure.* 1998;6(9):1153-67.
- 594 25. Buckle AM, Schreiber G, Fersht AR. Protein-protein recognition: crystal structural analysis of a  
595 barnase-barstar complex at 2.0-Å resolution. *Biochemistry.* 1994;33(30):8878-89. PubMed PMID:  
596 8043575.
- 597 26. Sogabe S, Stuart F, Henke C, Bridges A, Williams G, Birch A, et al. Neutralizing epitopes on the  
598 extracellular interferon  $\gamma$  receptor (IFN $\gamma$ R)  $\alpha$ -chain characterized by homolog scanning mutagenesis and X-  
599 ray crystal structure of the A6 Fab-IFN $\gamma$ R1-108 complex. Edited by R. Huber. *J Mol Biol.*  
600 1997;273(4):882-97.
- 601 27. Scheidig AJ, Hynes TR, Pelletier LA, Wells JA, Kossiakoff AA. Crystal structures of bovine  
602 chymotrypsin and trypsin complexed to the inhibitor domain of alzheimer's amyloid  $\beta$ -protein precursor  
603 (APPI) and basic pancreatic trypsin inhibitor (BPTI): Engineering of inhibitors with altered specificities.  
604 *Protein Sci.* 1997;6(9):1806-24.
- 605 28. Braden BC, Souchon H, Eiselé J-L, Bentley GA, Bhat TN, Navaza J, et al. Three-dimensional  
606 structures of the free and the antigen-complexed Fab from monoclonal anti-lysozyme antibody D44.1. *J*  
607 *Mol Biol.* 1994;243(4):767-81.
- 608 29. Hage T, Sebald W, Reinemer P. Crystal Structure of the Interleukin-4/Receptor  $\alpha$ 1 Chain  
609 Complex Reveals a Mosaic Binding Interface. *Cell.* 1999;97(2):271-81.
- 610 30. Bhat TN, Bentley GA, Boulot G, Greene MI, Tello D, Dall'Acqua W, et al. Bound water molecules  
611 and conformational stabilization help mediate an antigen-antibody association. *Proceedings of the National*  
612 *Academy of Sciences.* 1994;91(3):1089-93.
- 613 31. Lim D, Park HU, De Castro L, Kang SG, Lee HS, Jensen S, et al. Crystal structure and kinetic  
614 analysis of  $\beta$ -lactamase inhibitor protein-II in complex with TEM-1  $\beta$ -lactamase. *Nat Struct Biol.*  
615 2001;8(10):848-52.
- 616 32. Li S, Schmitz KR, Jeffrey PD, Wiltzius JJW, Kussie P, Ferguson KM. Structural basis for inhibition  
617 of the epidermal growth factor receptor by cetuximab. *Cancer Cell.* 2005;7(4):301-11.
- 618 33. Huang L, Hofer F, Martin GS, Kim S-H. Structural basis for the interaction of Ras with RaIGDS.  
619 *Nat Struct Biol.* 1998;5(6):422-6.
- 620 34. Prasad L, Waygood EB, Lee JS, Delbaere LTJ. The 2.5 Å resolution structure of the Jel42 Fab  
621 fragment/HPr complex. Edited by I. A. Wilson. *J Mol Biol.* 1998;280(5):829-45.
- 622 35. Bode W, Wei AZ, Huber R, Meyer E, Travis J, Neumann S. X-ray crystal structure of the complex  
623 of human leukocyte elastase (PMN elastase) and the third domain of the turkey ovomucoid inhibitor.  
624 *EMBO J.* 1986;5(10):2453-8. PubMed PMID: 3640709; PubMed Central PMCID: PMC1167139.
- 625 36. Padlan EA, Silvertown EW, Sheriff S, Cohen GH, Smith-Gill SJ, Davies DR. Structure of an  
626 antibody-antigen complex: crystal structure of the HyHEL-10 Fab-lysozyme complex. *Proc Natl Acad Sci*  
627 *U S A.* 1989;86(15):5938-42. PubMed PMID: 2762305; PubMed Central PMCID: PMC297746.
- 628 37. Meenan NAG, Sharma A, Fleishman SJ, MacDonald CJ, Morel B, Boetzel R, et al. The structural  
629 and energetic basis for high selectivity in a high-affinity protein-protein interaction. *Proceedings of the*  
630 *National Academy of Sciences.* 2010;107(22):10080-5.

- 631 38. Ultsch M, Bevers J, Nakamura G, Vandlen R, Kelley RF, Wu LC, et al. Structural Basis of  
632 Signaling Blockade by Anti-IL-13 Antibody Lebrikizumab. *J Mol Biol.* 2013;425(8):1330-9.
- 633 39. Guerois R, Nielsen JE, Serrano L. Predicting Changes in the Stability of Proteins and Protein  
634 Complexes: A Study of More Than 1000 Mutations. *Journal of Molecular Biology.* 2002;320(2):369-87.
- 635 40. Viricel C, de Givry S, Schiex T, Barbe S. Cost function network-based design of protein-protein  
636 interactions: predicting changes in binding affinity. *Bioinformatics.* 2018;34(15):2581-9.
- 637 41. Geng C, Vangone A, Folkers GE, Xue LC, Bonvin AMJJ. iSEE: Interface structure, evolution, and  
638 energy-based machine learning predictor of binding affinity changes upon mutations. *Proteins: Structure,*  
639 *Function, and Bioinformatics.* 2019;87(2):110-9.
- 640 42. Hurley B, O'Sullivan B, Allouche D, Katsirelos G, Schiex T, Zytnicki M, et al. Multi-language  
641 evaluation of exact solvers in graphical model discrete optimization. *Constraints.* 2016;21(3):413-34.
- 642 43. Alford RF, Leaver-Fay A, Jeliakov JR, O'Meara MJ, DiMaio FP, Park H, et al. The Rosetta All-  
643 Atom Energy Function for Macromolecular Modeling and Design. *J Chem Theory Comput.*  
644 2017;13(6):3031-48.
- 645 44. Park H, Bradley P, Greisen P, Liu Y, Mulligan VK, Kim DE, et al. Simultaneous Optimization of  
646 Biomolecular Energy Functions on Features from Small Molecules and Macromolecules. *J Chem Theory*  
647 *Comput.* 2016;12(12):6201-12.
- 648 45. Miller CR, Johnson EL, Burke AZ, Martin KP, Miura TA, Wichman HA, et al. Initiating a watch  
649 list for Ebola virus antibody escape mutations. *PeerJ.* 2016;4:e1674.
- 650 46. Patel JS, Quates CJ, Johnson EL, Ytreberg FM. Expanding the watch list for potential Ebola virus  
651 antibody escape mutations. *PLOS ONE.* 2019;14(3):e0211093.
- 652 47. Yang J, Naik N, Patel JS, Wylie CS, Gu W, Huang J, et al. Predicting the viability of beta-  
653 lactamase: How folding and binding free energies correlate with beta-lactamase fitness. *PloS one.*  
654 2020;15(5):e0233509.
- 655 48. Croux C, Dehon C. Influence functions of the Spearman and Kendall correlation measures.  
656 *Statistical Methods & Applications.* 2010;19(4):497-515.
- 657 49. Tien MZ, Meyer AG, Sydykova DK, Spielman SJ, Wilke CO. Maximum Allowed Solvent  
658 Accessibilities of Residues in Proteins. *PloS one.* 2013;8(11):e80635.
- 659 50. Kabsch W, Sander C. Dictionary of protein secondary structure: Pattern recognition of hydrogen-  
660 bonded and geometrical features. *Biopolymers.* 1983;22(12):2577-637.
- 661 51. Tribello GA, Bonomi M, Branduardi D, Camilloni C, Bussi G. PLUMED 2: New feathers for an  
662 old bird. *Computer Physics Communications.* 2014;185(2):604-13.
- 663 52. Sneath PHA. Relations between chemical structure and biological activity in peptides. *J Theor Biol.*  
664 1966;12(2):157-95.
- 665 53. van Zundert GCP, Rodrigues JPGLM, Trellet M, Schmitz C, Kastiris PL, Karaca E, et al. The  
666 HADDOCK2.2 Web Server: User-Friendly Integrative Modeling of Biomolecular Complexes. *J Mol Biol.*  
667 2016;428(4):720-5.
- 668 54. Dominguez C, Boelens R, Bonvin AMJJ. HADDOCK: A Protein-Protein Docking Approach  
669 Based on Biochemical or Biophysical Information. *J Am Chem Soc.* 2003;125(7):1731-7.
- 670 55. Wassenaar TA, van Dijk M, Loureiro-Ferreira N, van der Schot G, de Vries SJ, Schmitz C, et al.  
671 WeNMR: Structural Biology on the Grid. *Journal of Grid Computing.* 2012;10(4):743-67.
- 672 56. Patel JS, Ytreberg FM. Fast Calculation of Protein-Protein Binding Free Energies Using Umbrella  
673 Sampling with a Coarse-Grained Model. *J Chem Theory Comput.* 2018;14(2):991-7.
- 674  
675  
676  
677  
678  
679

680 **Supporting information captions**

681  
682 **S1 File. A word document with detailed protocols for calculating  $\Delta\Delta G$  values using each of**  
683 **the eight methods used in this study.**

684  
685 **S2 File. An in-house Python script that can be used to parse any of the parameters used in**  
686 **this study and provide tailored information.**

687  
688 **S3 File. A CSV file with full dataset used in this work and predicted  $\Delta\Delta G$  values for each**  
689 **mutation using eight methods.**

690  
691 **S1 Figure. Performance of each evaluated method for Ab and non-Ab complexes in**  
692 **predicting true  $\Delta\Delta G$  values ( $\rho_c$ ), linearly correlated  $\Delta\Delta G$  values ( $r$ ), and rank order ( $\rho$  and  $\tau$ )**  
693 **for a select subset of mutations that occur in beta sheet. The error for each method is reported**  
694 **under the correlation points.**

695  
696 **S2 Figure. Performance of each evaluated method for Ab and non-Ab complexes in**  
697 **predicting true  $\Delta\Delta G$  values ( $\rho_c$ ), linearly correlated  $\Delta\Delta G$  values ( $r$ ), and rank order ( $\rho$  and  $\tau$ )**  
698 **for a select subset of mutations that occur in alpha helix. The error for each method is reported**  
699 **under the correlation points.**

700  
701 **S3 Figure. Performance of each evaluated method for Ab and non-Ab complexes in**  
702 **predicting true  $\Delta\Delta G$  values ( $\rho_c$ ), linearly correlated  $\Delta\Delta G$  values ( $r$ ), and rank order ( $\rho$  and  $\tau$ )**  
703 **for a select subset of mutations with wild type amino acids that are either glycine or proline.**  
704 **The error for each method is reported under the correlation points.**

705  
706 **S4 Figure. Performance of each evaluated method for Ab and non-Ab complexes in**  
707 **predicting true  $\Delta\Delta G$  values ( $\rho_c$ ), linearly correlated  $\Delta\Delta G$  values ( $r$ ), and rank order ( $\rho$  and  $\tau$ )**  
708 **for a select subset of mutations with wild type amino acids that are neither glycine nor**  
709 **proline. The error for each method is reported under the correlation points.**

bioRxiv preprint doi: <https://doi.org/10.1101/2020.09.30.320069>; this version posted September 30, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

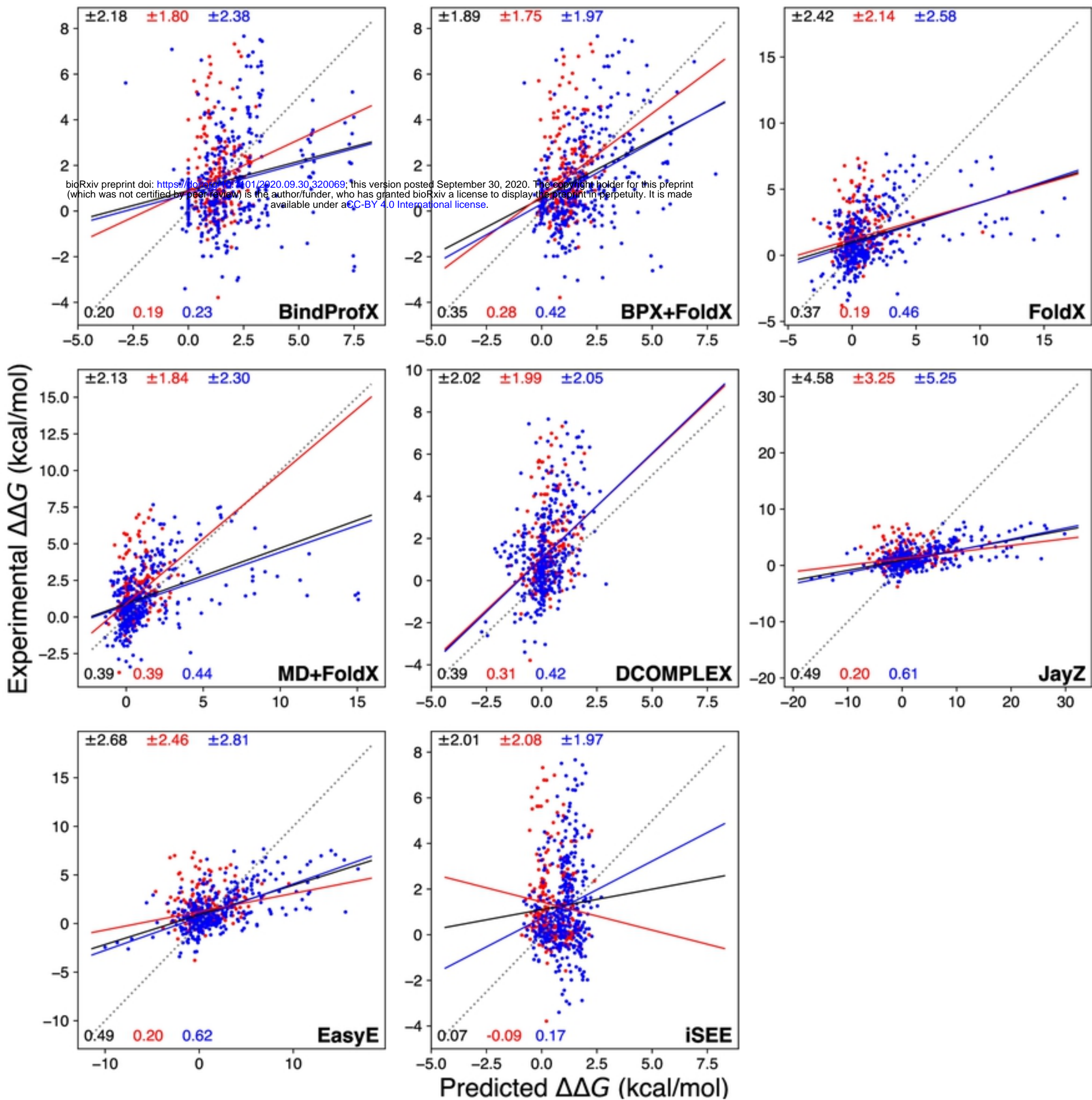


Figure 1

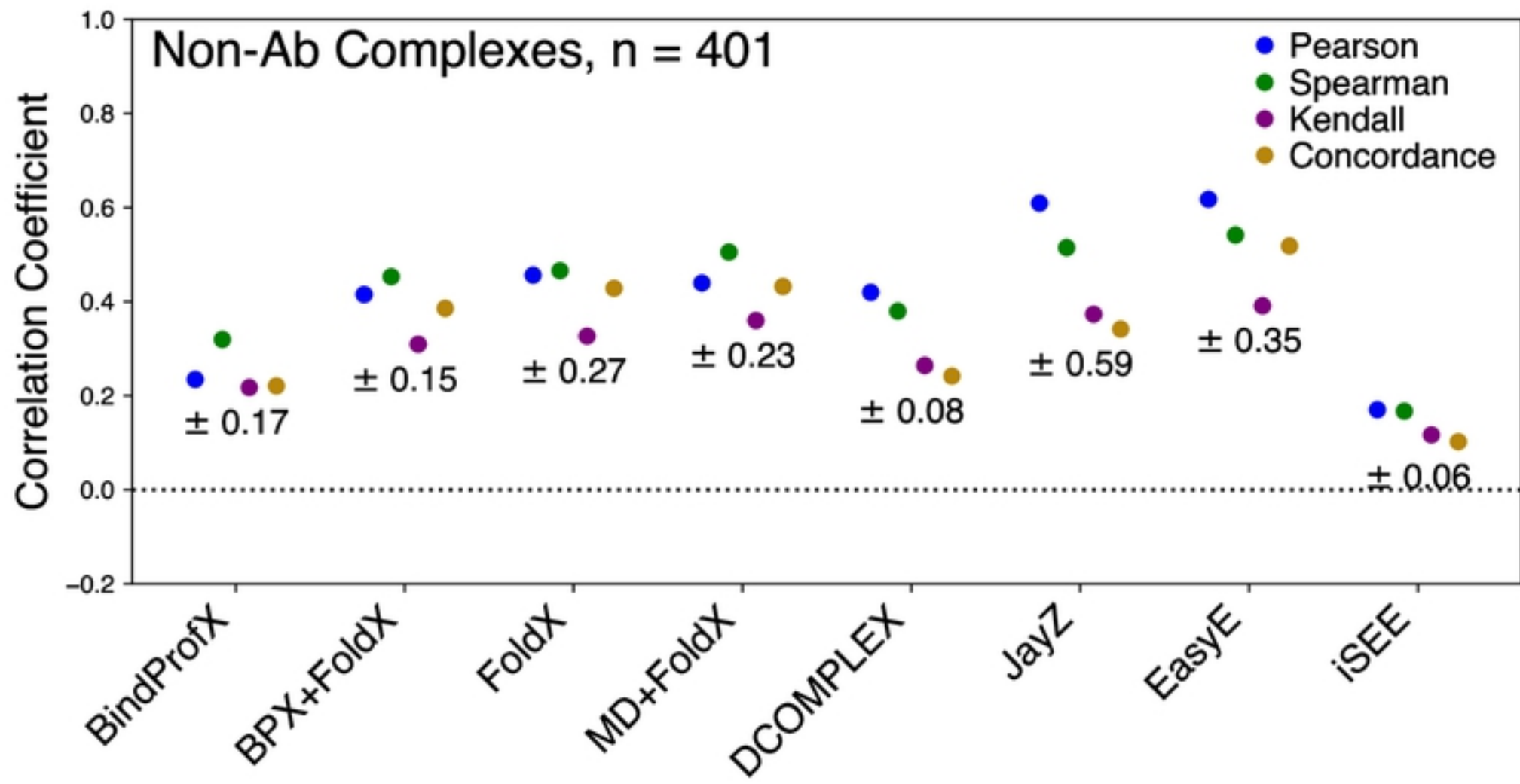


Figure 2

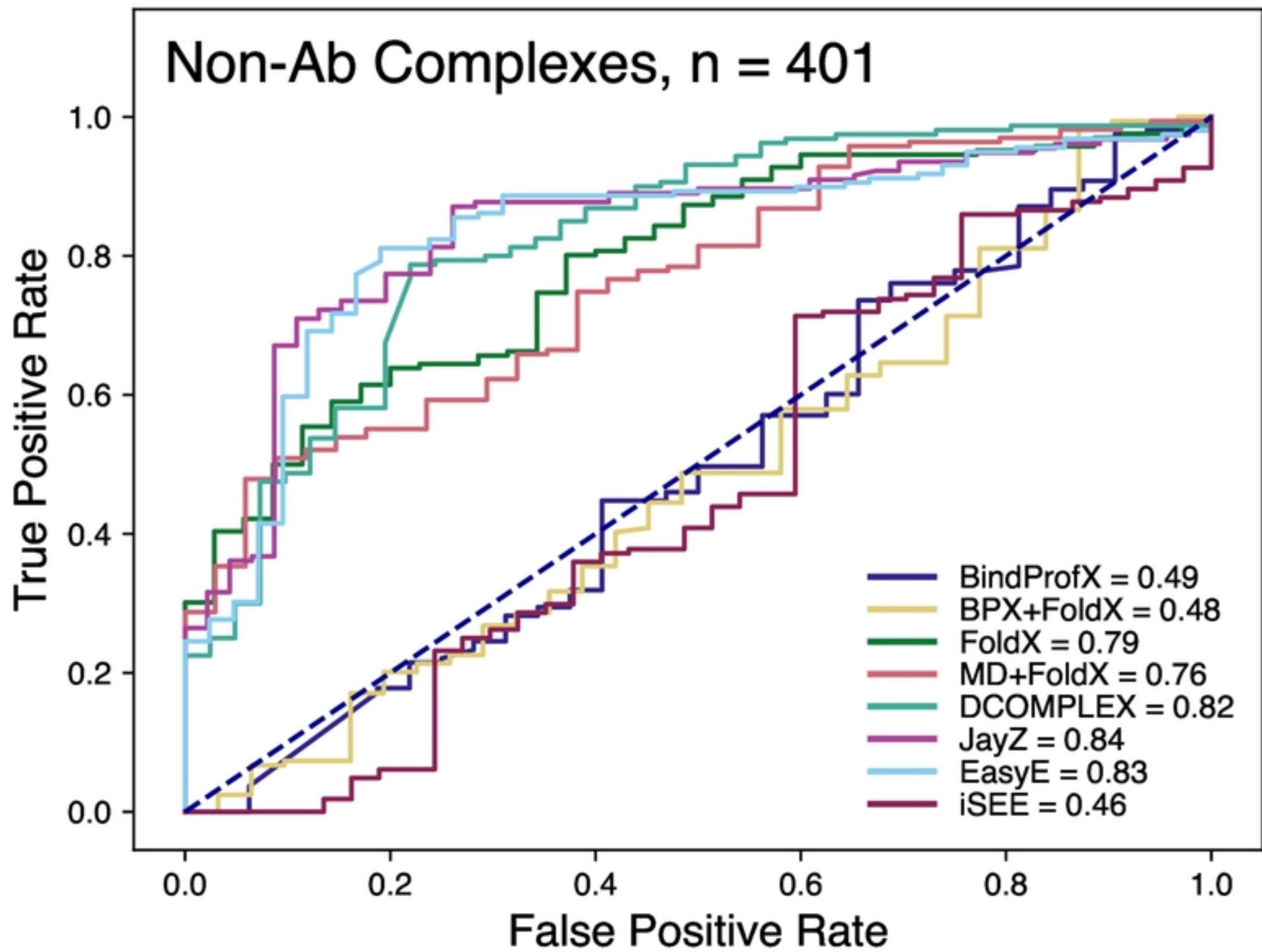


Figure 3

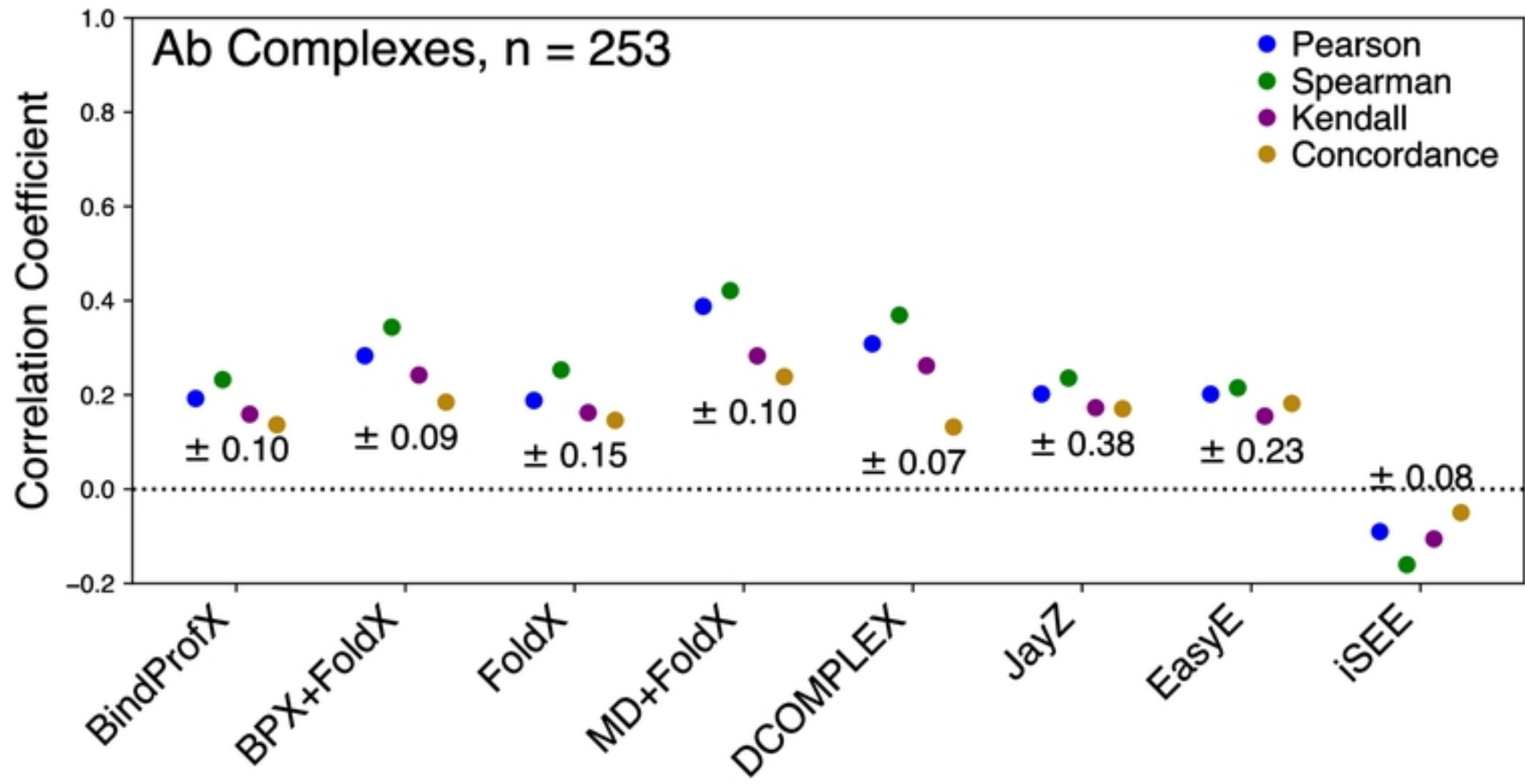


Figure 4



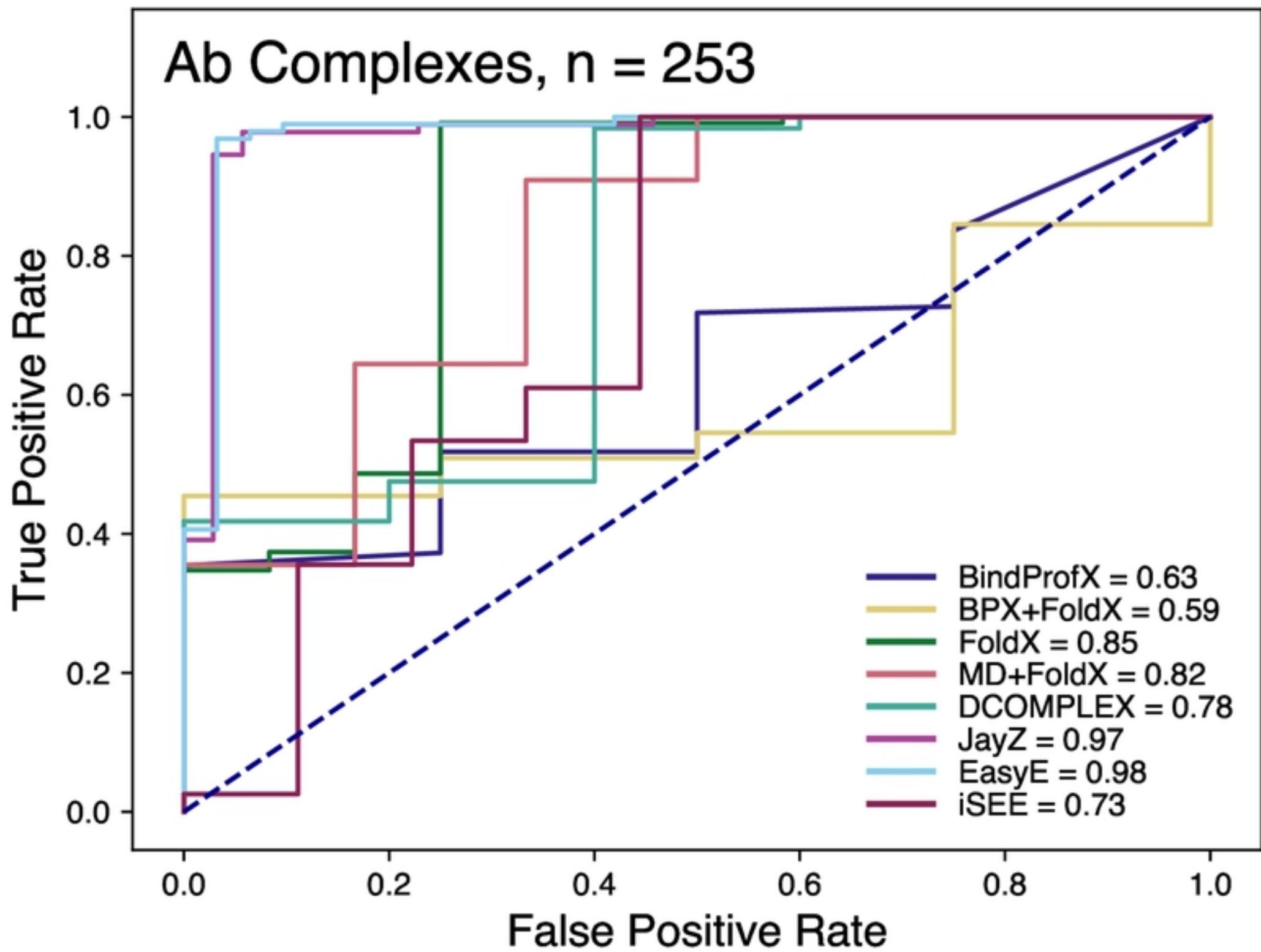


Figure 5