

000  
001  
002  
003        **Appendix — Enhancing Multimodal Protein Function**  
004        **Prediction through Dual-Branch Dynamic Selection with**  
005        **Reconstructive Pre-training**

009  
010        **1 Extended Experiments of Main Modules**

011        To further validate the rationality of the module design, we conducted an in-depth exploration of the  
012        two key components in DSRPGO: BInM and DSM.

013        **1.1 BInM Exploration**

014        To explore the feature interaction capability within BInM, we designed an experiment in which we  
015        separately removed the cross-attention module. DSRPGO-BInM-0 and DSRPGO-BInM-1 represent  
016        the removal of the bottom and top cross-attention modules from BInM.

017        Based on the results in Table 1, we find that removing either of the cross-attention modules leads to  
018        a decline in overall performance. This validates BInM’s ability to capture interactions.

019        **Table 1: Performance Comparison Under Different Interaction Settings.** The comparison results of  
020         $F_{max}$  score, micro-averaged AUPR (m-AUPR), macro-averaged AUPR (M-AUPR), F1-score (F1),  
021        and accuracy (ACC), for BPO, MFO, and CCO. The best results are highlighted in bold, and the  
022        sub-optimal results are underlined.

Method	F <sub>max</sub>			m-AUPR			M-AUPR			F1			ACC		
	BPO	MFO	CCO	BPO	MFO	CCO	BPO	MFO	CCO	BPO	MFO	CCO	BPO	MFO	CCO
DSRPGO	<b>0.458</b>	<b>0.254</b>	<b>0.452</b>	<b>0.330</b>	<b>0.166</b>	<b>0.371</b>	<b>0.182</b>	<b>0.114</b>	<b>0.239</b>	<b>0.272</b>	<b>0.241</b>	<b>0.357</b>	<b>0.346</b>	<b>0.124</b>	<b>0.262</b>
DSRPGO-BInM-0	0.433	0.189	0.351	0.314	0.113	0.282	0.175	0.106	0.192	0.271	0.179	0.306	0.296	0.086	0.166
DSRPGO-BInM-1	0.434	0.193	0.331	0.314	0.111	0.249	0.170	0.112	0.173	0.268	0.178	0.287	0.289	0.090	0.163

037  
038        **1.2 DSM Exploration**

039        To explore the impact of threshold  $t$  selection in the DSM module on model performance, we compare  
040        the model’s performance under different threshold values. Since there are six experts, the  
041        average selection probability for each expert is 1/6. Therefore, we set the threshold values around  
042        1/6. Notably, the DSRPGO model employs a learnable, adaptive threshold  $t$ .

043        Based on the experimental results in Table 2, we find that the model performs best when using a  
044        learnable, adaptive threshold  $t$ . This demonstrates the effectiveness of automatic threshold selection  
045        in the DSM module.

048  
049        **2 Extended Experiments of Pre-training Encoders**

050        In this work, the corresponding PSSI encoder and PSeI encoder are trained according to the spatial  
051        structure feature and sequence feature in the pre-training stage. To verify the effectiveness of the  
052        PSSI encoder and the PSeI encoder, the ablation experiment of the key module of the encoder is  
053        carried out in this section.

054  
 055 Table 2: Performance Comparison of Different Threshold  $t$  in DSM Module. The comparison  
 056 results of  $F_{max}$  score, micro-averaged AUPR (m-AUPR), macro-averaged AUPR (M-AUPR), F1-  
 057 score (F1), and accuracy (ACC), for BPO, MFO, and CCO. The best results are highlighted in bold.  
 058

Threshold $t$	F <sub>max</sub>			m-AUPR			M-AUPR			F1			ACC		
	BPO	MFO	CCO	BPO	MFO	CCO	BPO	MFO	CCO	BPO	MFO	CCO	BPO	MFO	CCO
auto (ours)	<b>0.458</b>	<b>0.254</b>	<b>0.452</b>	<b>0.330</b>	<b>0.166</b>	<b>0.371</b>	<b>0.182</b>	<b>0.114</b>	<b>0.239</b>	<b>0.272</b>	<b>0.241</b>	<b>0.357</b>	<b>0.346</b>	<b>0.124</b>	<b>0.262</b>
0.1	0.430	0.202	0.428	0.326	0.122	0.314	0.103	0.103	0.221	0.266	0.189	0.325	0.341	0.097	0.252
0.14	0.429	0.201	0.386	0.325	0.123	0.297	0.109	0.109	0.216	0.269	0.187	0.325	0.330	0.099	0.252
0.1616	0.433	0.204	0.427	0.327	0.123	0.286	0.110	0.110	0.212	0.264	0.189	0.333	0.335	0.097	0.261
0.18	0.438	0.205	0.416	0.328	0.123	0.279	0.101	0.101	0.225	0.266	0.190	0.300	0.335	0.102	0.227
0.2	0.429	0.204	0.380	0.326	0.123	0.261	0.107	0.107	0.215	0.269	0.189	0.317	0.346	0.102	0.227
0.3	0.426	0.206	0.411	0.324	0.122	0.271	0.105	0.105	0.224	0.272	0.189	0.333	0.330	0.100	0.227
0.4	0.430	0.205	0.405	0.320	0.121	0.307	0.113	0.113	0.214	0.272	0.193	0.329	0.308	0.100	0.252

061  
 062  
 063  
 064  
 065  
 066  
 067 Table 3: Ablation experiments of key components in the pre-trained encoder. The comparison  
 068 results of  $F_{max}$  score, micro-averaged AUPR (m-AUPR), macro-averaged AUPR (M-AUPR), F1-  
 069 score (F1), and accuracy (ACC), for BPO, MFO, and CCO. The Cost Time (ms) represents the time  
 070 it takes for the model to test each protein sample. The best results are highlighted in bold.  
 071

Method	F <sub>max</sub>			m-AUPR			Cost Time (ms)
	BPO	MFO	CCO	BPO	MFO	CCO	
Spatial-BiMamba Block	0.370	<b>0.240</b>	<b>0.419</b>	0.244	<b>0.156</b>	<b>0.371</b>	<b>0.539</b>
Spatial-Multihead Attention	<b>0.430</b>	0.223	0.350	<b>0.291</b>	0.154	0.313	5.661
Sequence-BiMamba Block	0.290	<b>0.237</b>	0.308	0.185	<b>0.157</b>	0.230	<b>0.280</b>
Sequence-Multihead Attention	<b>0.329</b>	0.219	<b>0.345</b>	<b>0.199</b>	0.148	<b>0.248</b>	0.450

072  
 073 In this experiment, the pre-training framework is the same as stage 1 of DSRPGO. Here, BiMamba  
 074 Block in the pre-trained encoder can be replaced by Multihead Attention Block Wu et al. (2023),  
 075 and vice versa. Then, we use the pre-trained encoder for feature extraction in the fine-tuning task.  
 076 During fine-tuning, we only use an MLP for classification.

077  
 078 As shown in Table 3, the model names are formed by combining the feature and framework compo-  
 079 nents with a hyphen. This indicates that the results are obtained by pre-training the feature with the  
 080 framework component and then fine-tuning it.

081 We model the spatial structure features and sequence features using encoders with different compo-  
 082 nents. Experimental results in Table 3 show that the encoder using BiMamba Block as a component  
 083 for modeling spatial structure features shows significant advantages in MFO and CCO. When mod-  
 084 eling sequence features, the encoder utilizing Multihead Attention as a component demonstrates  
 085 considerable advantages in BPO and CCO. Additionally, we observe that BiMamba Block is signifi-  
 086 cantly more efficient in processing spatial structure information, requiring much less time to predict  
 087 a protein’s function compared to Multihead Attention. Considering both inference efficiency and  
 088 performance, we decide to use an encoder with BiMamba Block as the component for spatial struc-  
 089 ture features modeling. For sequence information, while BiMamba Block achieves better inference  
 090 efficiency, we prioritize model performance and opt for Multihead Attention to model sequence  
 091 features.

### 3 Impact of Sequence Similarity

101 To ensure the validity of our experimental design and avoid potential data leakage, we analyze the  
 102 sequence similarity between the test set and the combined training and validation sets for BPO,  
 103 MFO, and CCO. We calculate the similarity for each of these sets and categorize the results into  
 104 different similarity ranges.

105 Figure 1 shows the number of proteins in the test set within each similarity range. We observe  
 106 that the majority of proteins in our test set exhibit an average sequence similarity of less than 50%  
 107 with the proteins in the combined training and validation sets. Only a few proteins have an average

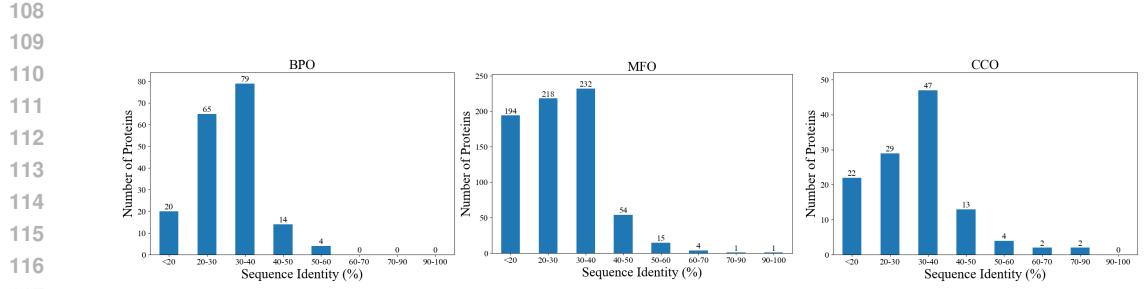


Figure 1: Distribution of sequence identity across proteins in the test dataset. The x-axis represents the sequence similarity ranges. The y-axis indicates the number of proteins within each range. Each bar denotes the number of proteins in the test set that fall within the corresponding similarity range.

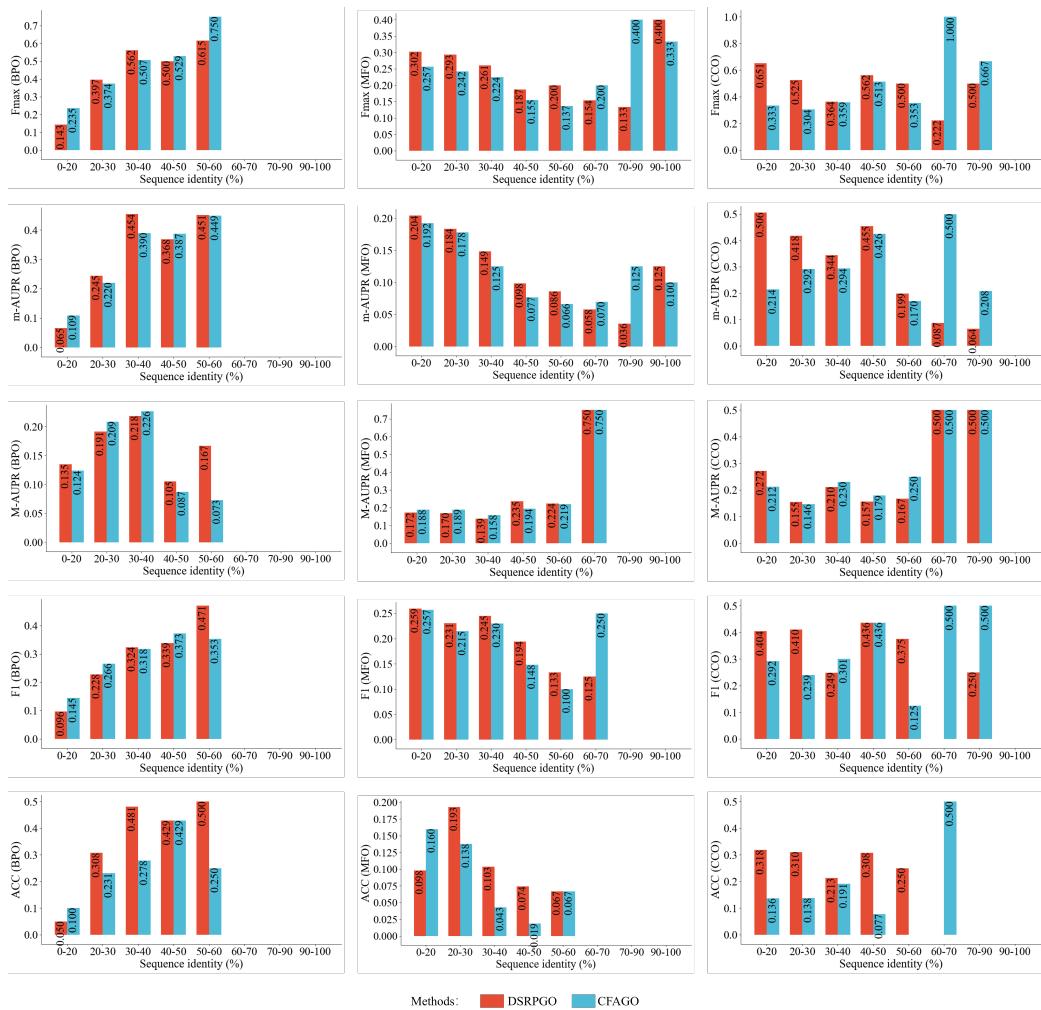


Figure 2: Performance of the different methods in Different Sequence Similarity Ranges. Figures show the results of F1-score (F1), accuracy (ACC), F-max score, macro-averaged AUPR (M-AUPR), and micro-averaged AUPR (m-AUPR) for BPO, MFO and CCO, respectively. The red pillar represents the results of the DSRPGO model and the blue pillar shows the results of the CFAGO model.

162  
163 Table 4: Performance of the different methods in Different Sequence Similarity Ranges. The  
164 comparison results of  $F_{max}$  score, micro-averaged AUPR (m-AUPR), macro-averaged AUPR (M-  
165 AUPR), F1-score (F1), and accuracy (ACC), for BPO, MFO, and CCO. The Range represents the  
proteins that have average similarity within a specific range.

Method	Range	$F_{max}$			m-AUPR			M-AUPR			F1			Acc		
		BPO	MFO	CCO	BPO	MFO	CCO	BPO	MFO	CCO	BPO	MFO	CCO	BPO	MFO	CCO
CFAGO	0-100	0.439	0.236	0.366	0.328	0.159	0.337	0.188	0.138	0.210	0.283	0.234	0.314	0.338	0.100	0.210
	0-20	0.235	0.257	0.333	0.109	0.192	0.214	0.124	0.188	0.212	0.145	0.257	0.292	0.100	0.160	0.136
	20-30	0.374	0.242	0.304	0.220	0.178	0.292	0.209	0.189	0.146	0.266	0.215	0.239	0.231	0.138	0.138
	30-40	0.507	0.224	0.359	0.390	0.125	0.294	0.226	0.158	0.230	0.318	0.230	0.301	0.278	0.043	0.191
	40-50	0.529	0.155	0.513	0.387	0.077	0.426	0.087	0.194	0.179	0.373	0.148	0.436	0.429	0.019	0.077
	50-60	0.750	0.137	0.353	0.449	0.066	0.170	0.073	0.219	0.250	0.353	0.100	0.125	0.250	0.067	0.000
	60-70	-	0.200	1.000	-	0.070	0.500	-	0.750	0.500	-	0.250	0.500	-	0.000	0.500
	70-90	-	0.400	0.667	-	0.125	0.208	-	0.000	0.500	-	0.000	0.500	-	0.000	0.000
	90-100	-	0.333	-	-	0.100	-	-	0.000	-	-	0.000	-	-	0.000	-
DSRPGO (Ours)	0-100	0.458	0.330	0.452	0.330	0.166	0.371	0.182	0.114	0.239	0.272	0.241	0.357	0.346	0.124	0.262
	0-20	0.143	0.302	0.651	0.065	0.204	0.506	0.135	0.172	0.272	0.096	0.259	0.404	0.05	0.098	0.318
	20-30	0.397	0.293	0.525	0.245	0.184	0.418	0.191	0.17	0.155	0.228	0.231	0.41	0.308	0.193	0.31
	30-40	0.562	0.261	0.364	0.454	0.149	0.344	0.218	0.139	0.21	0.324	0.245	0.249	0.481	0.103	0.213
	40-50	0.5	0.187	0.562	0.368	0.098	0.455	0.105	0.235	0.157	0.339	0.194	0.436	0.429	0.074	0.308
	50-60	0.615	0.2	0.5	0.451	0.086	0.199	0.167	0.224	0.167	0.471	0.133	0.375	0.5	0.067	0.25
	60-70	-	0.154	0.222	-	0.058	0.087	-	0.75	0.5	-	0.125	0	-	0	0
	70-90	-	0.133	0.5	-	0.036	0.064	-	0	0.5	-	0	0.25	-	0	0
	90-100	-	0.4	-	-	0.125	-	-	0	-	-	0	-	-	0	-

181  
182 Table 5: Comparison with Structure-based and PLM-based methods. The comparison results of  
183  $F_{max}$  score, and micro-averaged AUPR (m-AUPR) for BPO, MFO, and CCO. The best results are  
184 highlighted in bold.

Method	Focus	$F_{max}$			m-AUPR		
		BPO	MFO	CCO	BPO	MFO	CCO
DeepFRI	Structure based	0.362	<b>0.461</b>	0.385	0.308	<b>0.382</b>	0.360
PredGO	Structure + PLM based	0.108	0.455	0.252	0.058	0.254	0.183
DSRPGO (Ours)	Multi-modal based	<b>0.440</b>	0.282	<b>0.421</b>	<b>0.332</b>	0.172	<b>0.392</b>

191 similarity greater than 70%. Based on these results, we conclude that the time-based split used in  
192 the CAFA challenge is reasonable and does not introduce significant sequence similarity overlap  
193 between our training and test sets.

194 In the comparative experiments described in Section 3.2, BLAST, which primarily relies on se-  
195 quence similarity, performs poorly. This further supports the notion that the sequence similarity  
196 between the test set and the combined training and validation sets is relatively low.

197 Additionally, we explore the model’s performance in predicting protein functions within different  
198 similarity ranges. We divide the test set into 9 similarity intervals.

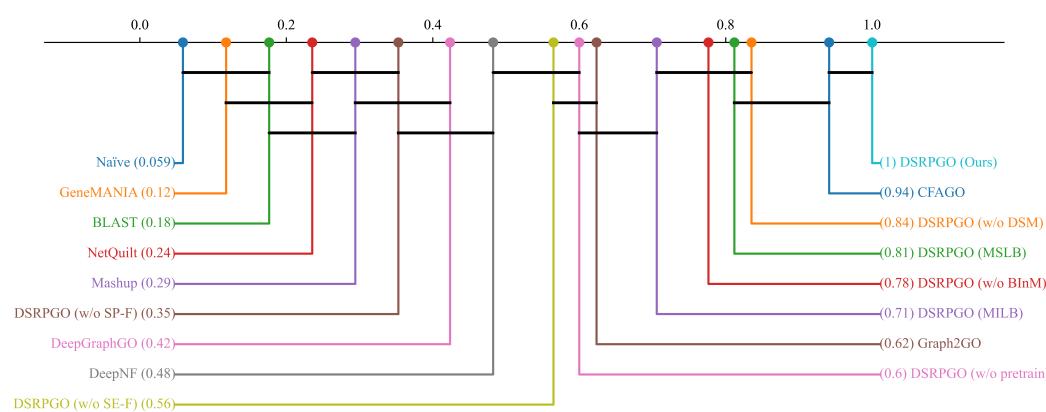
199 In Figure 2, we present the performance of DSRPGO and CFAGO across different protein similarity  
200 ranges. Specifically, range (n-m) represents the average similarity of proteins in the range of n-m.  
201 For both BPO and CCO, our method demonstrates significantly better performance than CFAGO,  
202 particularly for proteins with low similarity. And the detailed results are shown in Table 4. The  
203 symbol “-” indicates no proteins in that interval. Due to the scarcity of proteins with similarity  
204 above 50%, we do not provide further analysis for this group. For proteins with similarity below  
205 50%, the model performs best in predicting proteins in the similarity range of 30 to 40 for BPO. For  
206 MFO and CCO, the model performs best for proteins in the similarity range of 0 to 30. These results  
207 show that our proposed DSRPGO model is stable across different sequence similarities and is not  
208 reliant on the sequence similarity of the data.

## 210 4 Comparison with Structure-based and PLM-based Methods

211 To ensure a fair comparison, we conduct an additional experiment, evaluating the performance of  
212 the structure-based method DeepFRI (Gligorijević et al., 2021), and PredGO (Zheng et al., 2023)  
213 method using a protein language model (PLM). The results shown in Table 5 demonstrate that our  
214 approach outperforms both methods in BPO and CCO aspects.

216  
217 Table 6: Performance Comparison Under Different Learning Rate Settings. The comparison results  
218 of  $F_{max}$  score, micro-averaged AUPR (m-AUPR), macro-averaged AUPR (M-AUPR), F1-score  
219 (F1), and accuracy (ACC), for BPO, MFO, and CCO. The best results are highlighted in bold.

LR	F <sub>max</sub>			m-AUPR			M-AUPR			F1			ACC		
	BPO	MFO	CCO	BPO	MFO	CCO	BPO	MFO	CCO	BPO	MFO	CCO	BPO	MFO	CCO
1.00E-03 (ours)	<b>0.458</b>	0.254	<b>0.452</b>	<b>0.330</b>	<b>0.166</b>	<b>0.371</b>	0.182	0.114	0.239	0.272	0.241	<b>0.357</b>	<b>0.346</b>	<b>0.124</b>	<b>0.262</b>
1.00E-04	0.352	<b>0.265</b>	0.387	0.254	0.160	0.291	<b>0.195</b>	0.114	0.242	0.261	<b>0.261</b>	0.333	0.220	0.071	0.261
1.00E-05	0.063	0.146	0.291	0.034	0.055	0.262	0.101	0.064	0.231	0.035	0.009	0.202	0.000	0.000	0.134
2.00E-03	0.451	0.242	0.423	0.310	0.145	0.361	0.169	0.113	0.225	<b>0.280</b>	0.246	0.325	0.302	0.078	0.185
2.00E-04	0.422	0.231	0.435	0.307	0.148	0.324	0.187	<b>0.123</b>	0.225	0.269	0.234	0.337	0.264	0.085	0.244
2.00E-05	0.181	0.166	0.311	0.082	0.073	0.291	0.125	0.097	<b>0.264</b>	0.132	0.089	0.239	0.011	0.001	0.151



227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244 Figure 3: Critical Difference Diagrams.

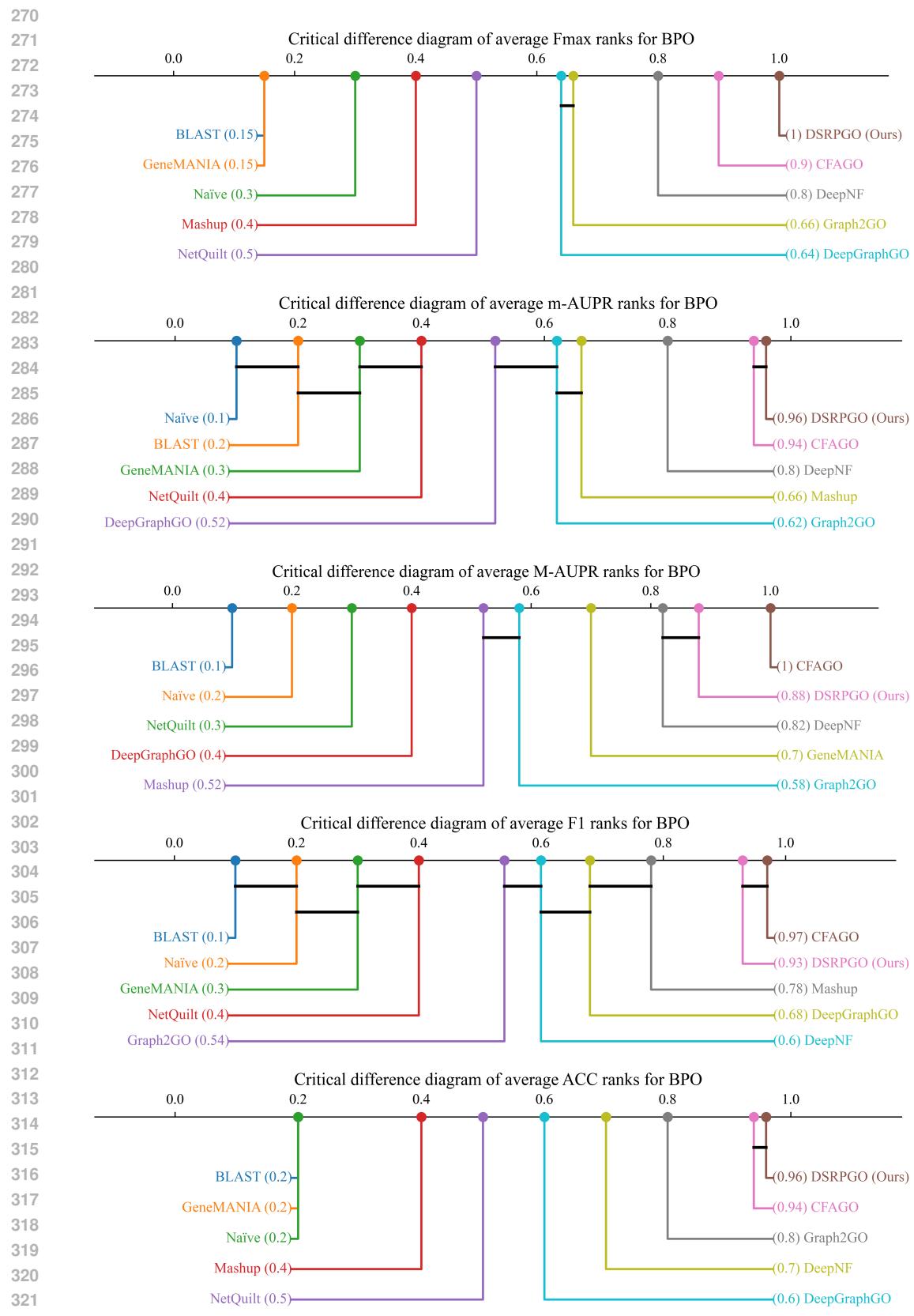
## 5 Impact of Learning Rate

246 To investigate the impact of the learning rate (LR) on model performance, we conduct experiments  
247 with various learning rates ranging from 1e-5 to 2e-3. All other hyperparameters, including the  
248 optimizer, batch size, and training epochs, are kept constant to isolate the effect of the learning rate.  
249 From Table 6, we can see that the model performs best when the learning rate is set to 1e-3. This  
250 may be because a lower learning rate leads to underfitting of the model.

## 6 Critical Difference Diagrams for Statistical Comparison

254 To assess the significance of the results and compare the performance of different approaches, we  
255 use critical difference diagrams (CDDs) (Benavoli et al., 2016). These diagrams, as described in  
256 scikit-posthocs documentation, are particularly useful in visualizing whether differences between  
257 approaches are statistically significant.

258 The CDD, as shown in Figure 3, presents the overall comparison of all models across various as-  
259 pects and metrics. To provide a more detailed analysis, we include additional CDDs to compare  
260 the models' performance across different aspects, metrics, and experiments. The CDDs used for  
261 the comparative experiments are shown in Figures 4, 5 and 6. The CDDs used for the module ab-  
262 lation experiments are shown in Figures 7, 8 and 9. By examining the CDDs, we can observe that  
263 DSRPGO ranks first in almost all cases, indicating that our model demonstrates significant advan-  
264 tages in single-species protein function prediction.



323                   Figure 4: Critical Difference Diagram for Comparative Experiments on BPO.

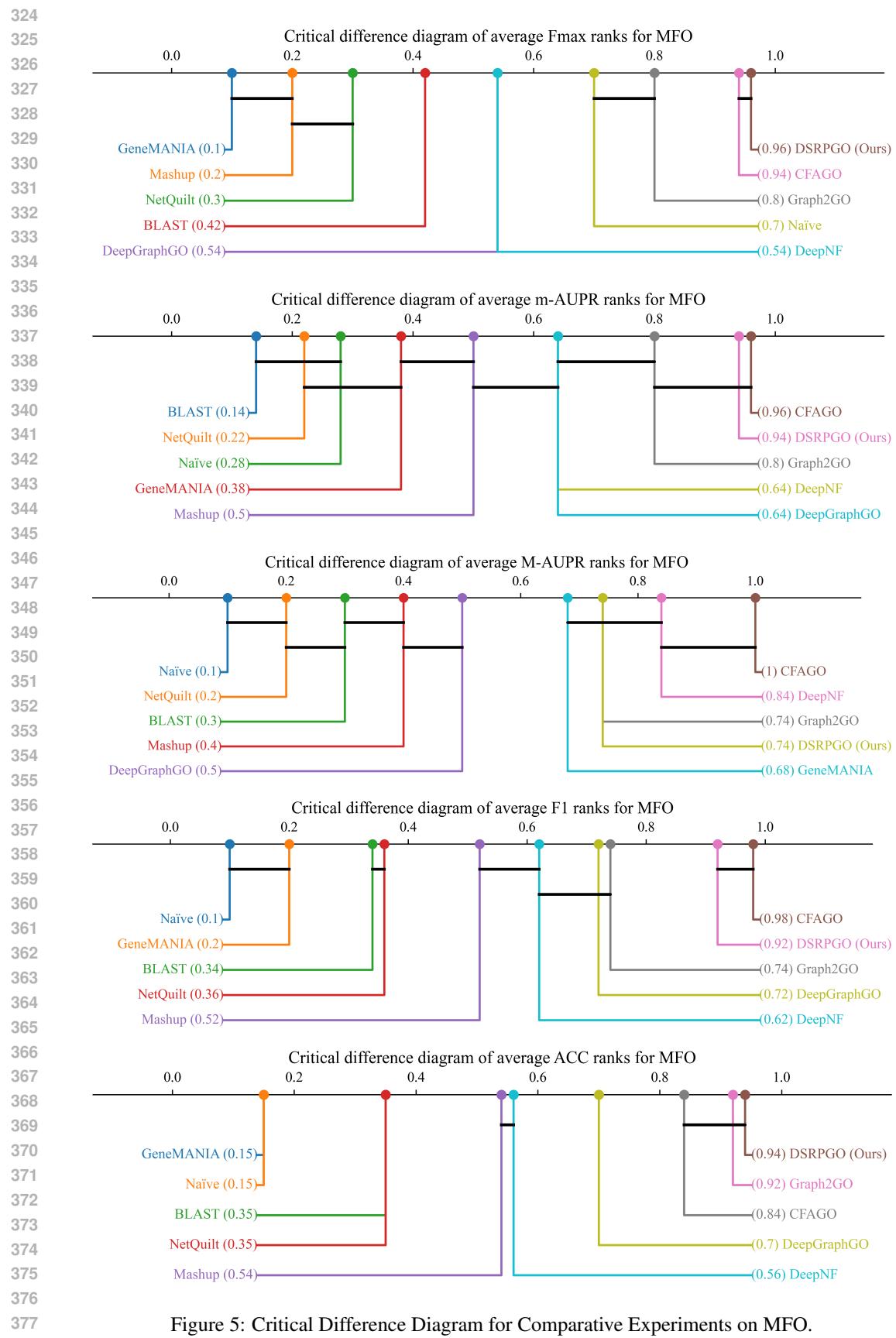


Figure 5: Critical Difference Diagram for Comparative Experiments on MFO.

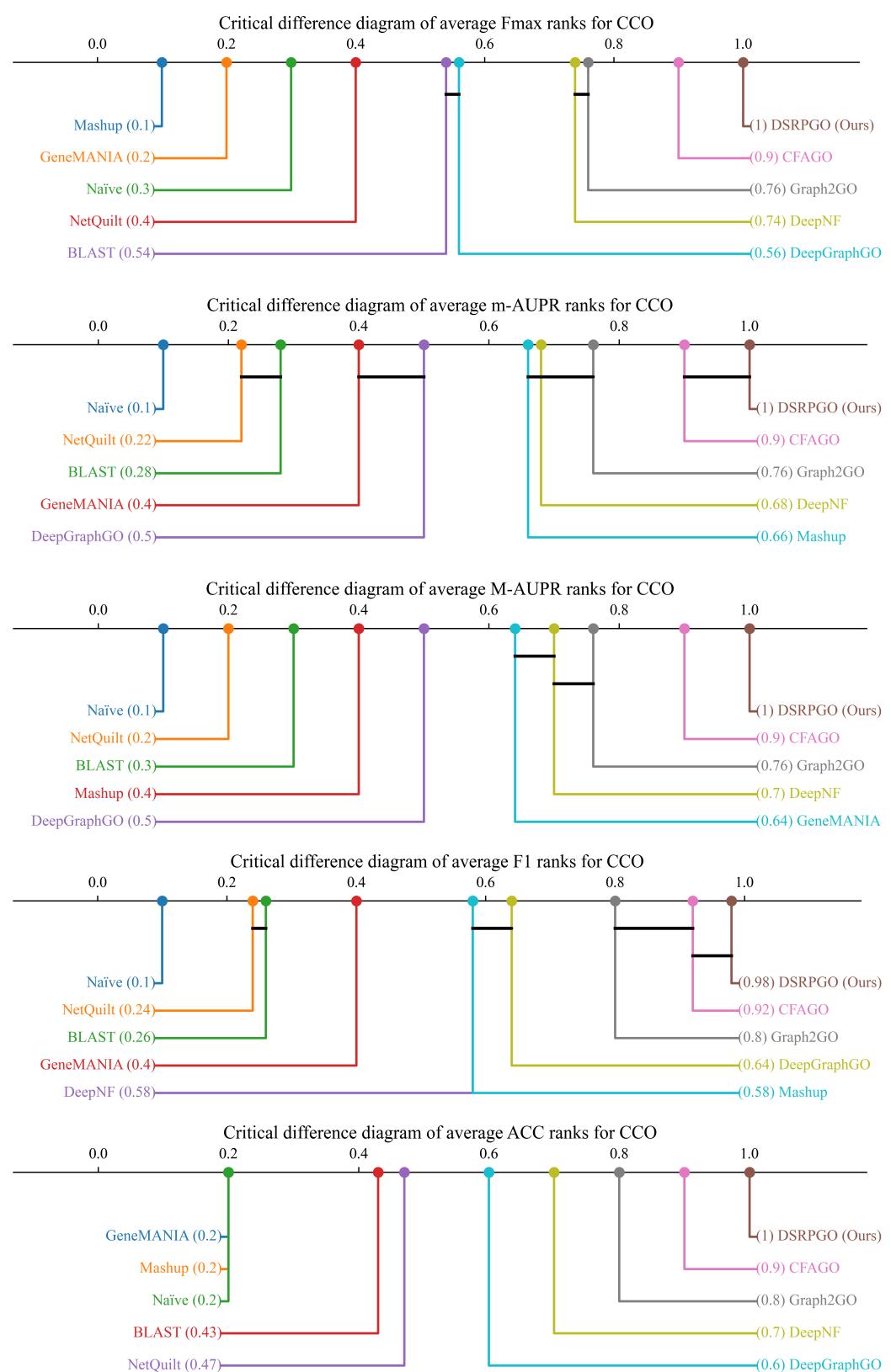
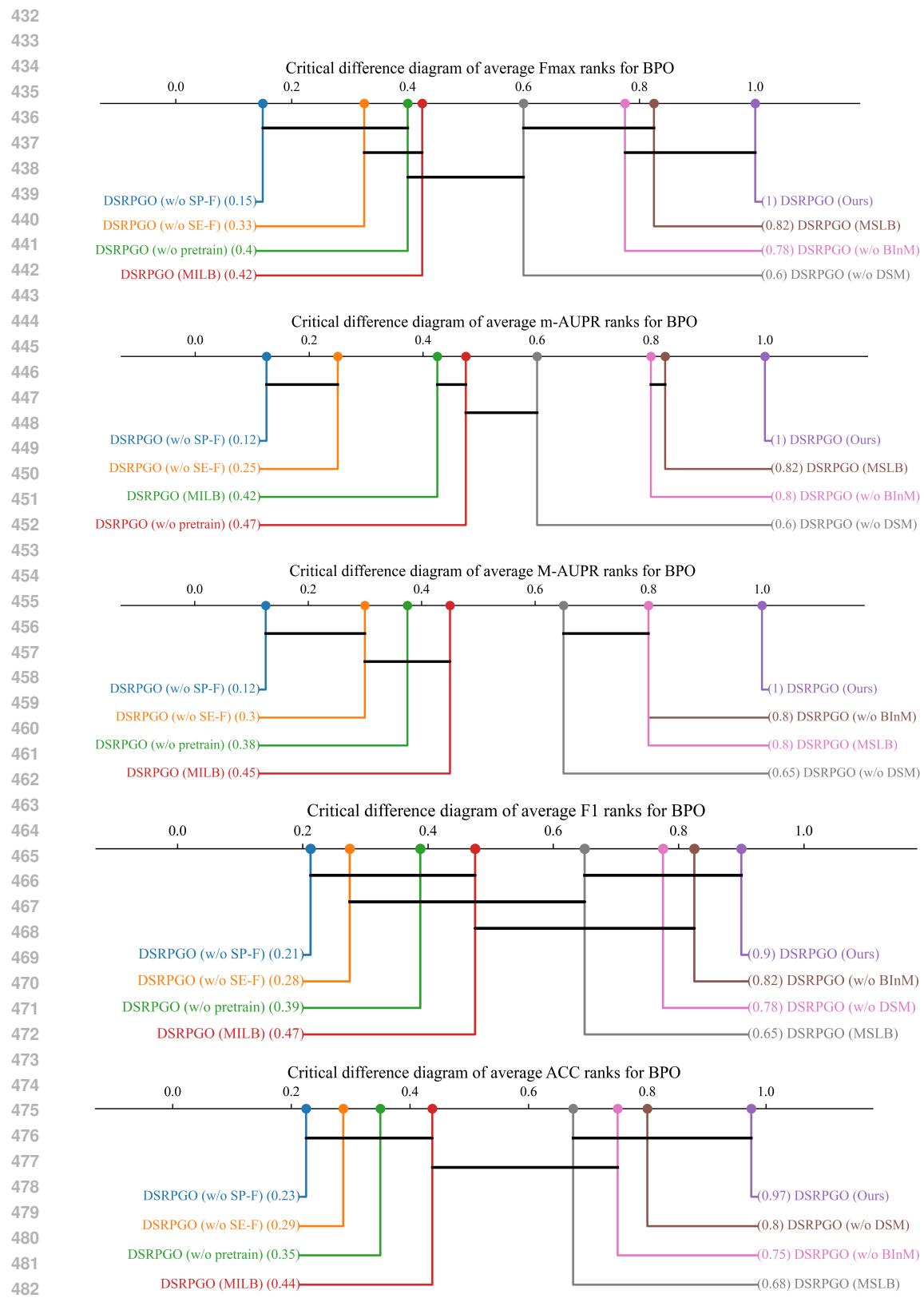


Figure 6: Critical Difference Diagram for Comparative Experiments on CCO.



483                    Figure 7: Critical Difference Diagram for Module Ablation on BPO.  
484  
485

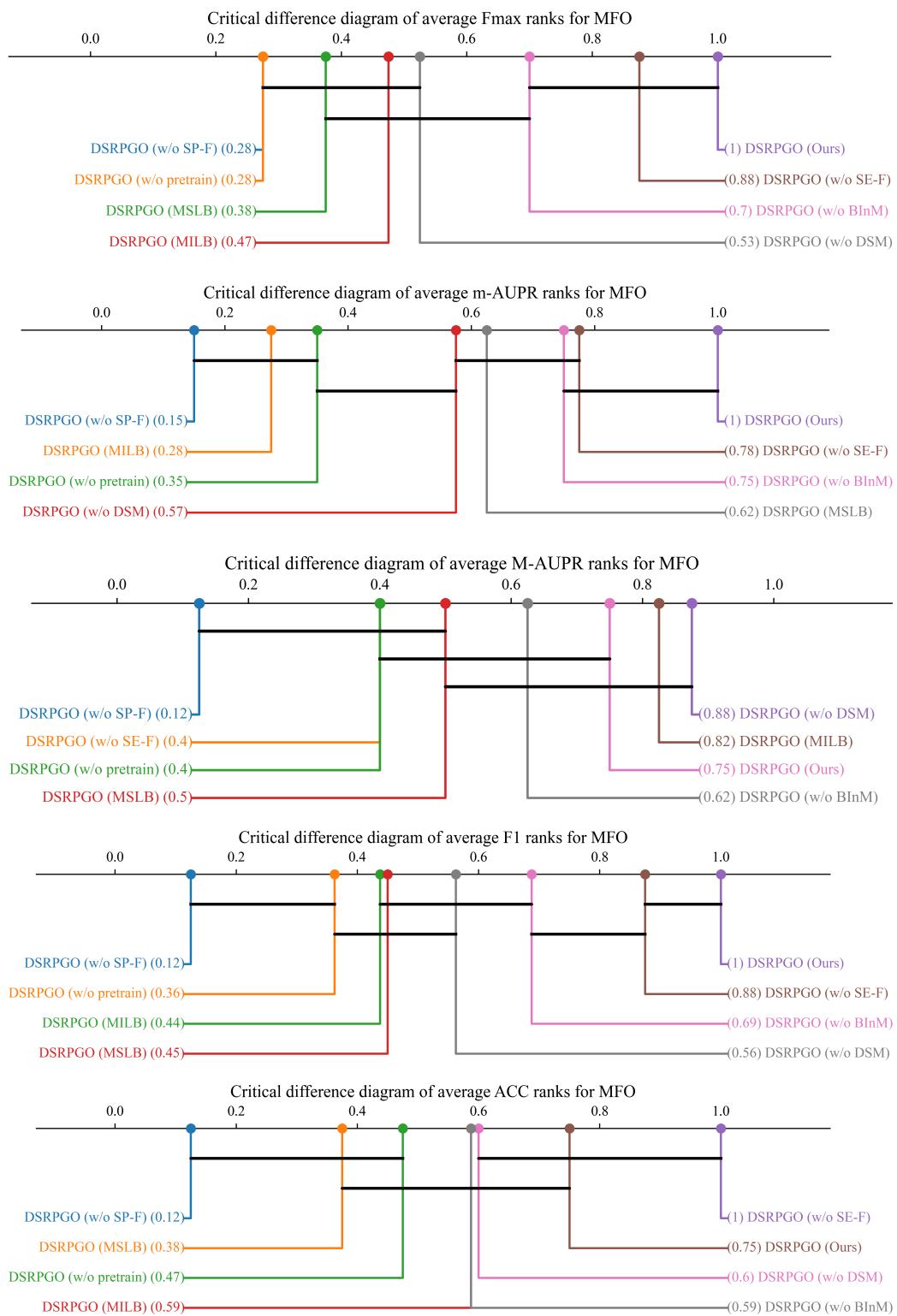
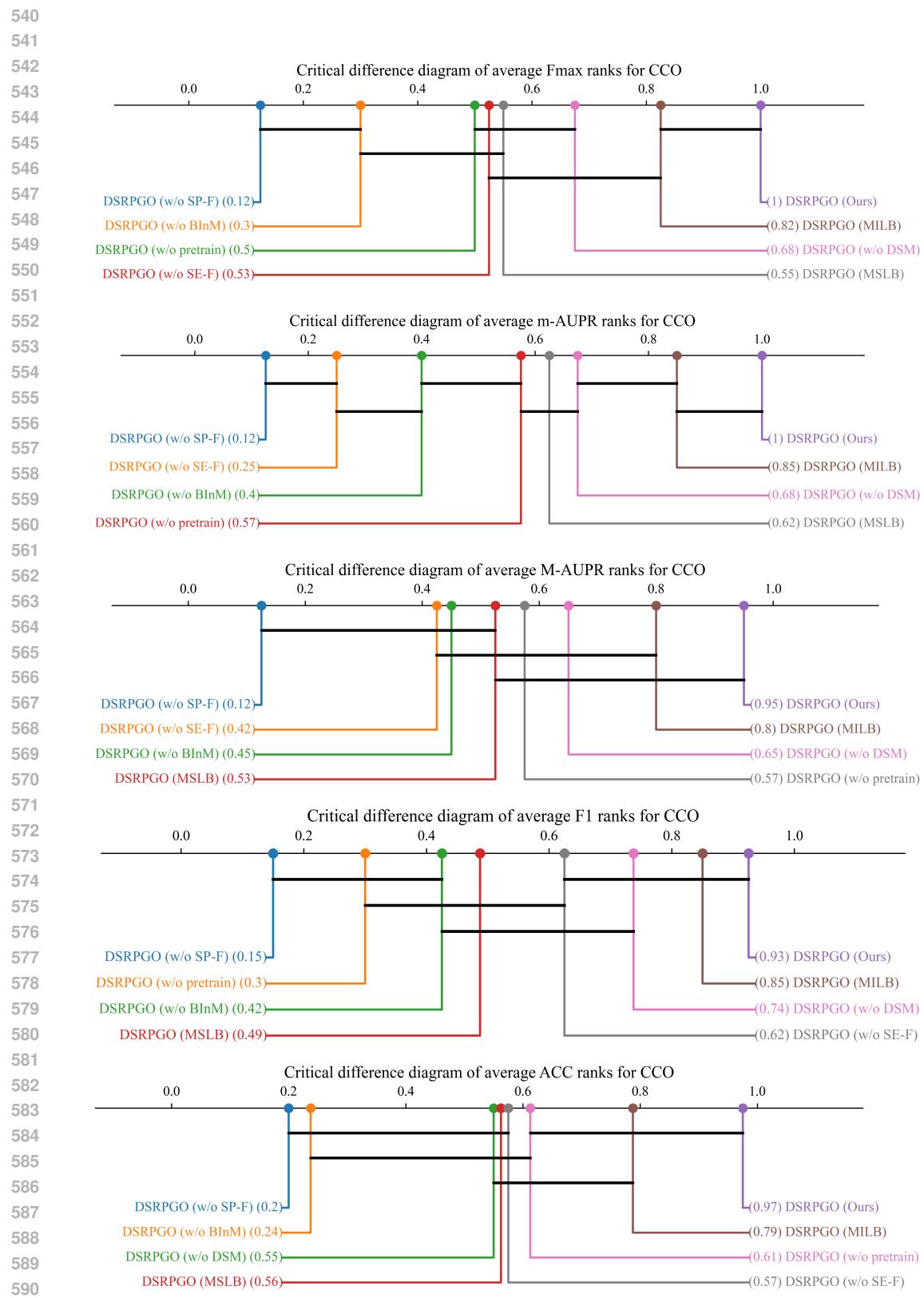


Figure 8: Critical Difference Diagram for Module Ablation on MFO.



591 Figure 9: Critical Difference Diagram for Module Ablation on CCO.  
592  
593

594      **References**  
595

596      Alessio Benavoli, Giorgio Corani, and Francesca Mangili. Should we really use post-hoc tests based  
597      on mean-ranks? *The Journal of Machine Learning Research*, 17(1):152–161, 2016.

598      Vladimir Gligorijević, P Douglas Renfrew, Tomasz Kosciolek, Julia Koehler Leman, Daniel Beren-  
599      berg, Tommi Vatanen, Chris Chandler, Bryn C Taylor, Ian M Fisk, Hera Vlamakis, et al. Structure-  
600      based protein function prediction using graph convolutional networks. *Nature communications*,  
601      12(1):3168, 2021.

602      Zhourun Wu, Mingyue Guo, Xiaopeng Jin, Junjie Chen, and Bin Liu. Cfago: cross-fusion of net-  
603      work and attributes based on attention mechanism for protein function prediction. *Bioinformatics*,  
604      39(3):btad123, 2023.  
605

606      Rongtao Zheng, Zhijian Huang, and Lei Deng. Large-scale predicting protein functions through  
607      heterogeneous feature fusion. *Briefings in Bioinformatics*, 24(4):bbad243, 2023.

608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647