

**School of Biological Sciences****ASSESSMENT COVER SHEET AND TEMPLATE****Section A – to be completed by the student**

Student Number	2456386		
Programme	<b>Msc Bioinformatics</b>		
Unit Name	<b>Bioinformatics Research Project</b>	Unit Code:	<b>BIOLM0034</b>
Assessment name	<b>Research Project Report</b>		
Word Count	<b>5921</b>		
Do you give permission for your work to be used anonymously in examples given to students in the future? Yes			

**By submitting this assignment cover sheet, I confirm that I understand and agree with the following statements:**

- 'I have not committed plagiarism, cheated or otherwise committed academic misconduct as defined in the University's Assessment Regulations (available at <https://www.bristol.ac.uk/media-library/sites/academic-quality/documents/taught-code/annexes/university-examination-regulations.pdf>)'
- 'I have not submitted this piece, in part or in its entirety, for assessment in another unit assignment (including at other institutions) as outlined in section 4 of the University's Assessment regulations (available at <https://www.bristol.ac.uk/media-library/sites/academic-quality/documents/taught-code/annexes/university-examination-regulations.pdf>)'
- 'I understand that this piece will be scrutinised by anti-plagiarism software and that I may incur penalties if I am found to have committed plagiarism, as outlined in sections 3 of the University's Examination Regulations (available at <https://www.bristol.ac.uk/media-library/sites/academic-quality/documents/taught-code/annexes/university-examination-regulations.pdf>)'

## CHALLENGES AND STRATEGIES IN FINDING DIELS- ALDER ENZYMES: A COMPARATIVE STUDY OF STRUCTURAL AND SEQUENCE-BASED METHODS



Nicole Alejandra Morveli Flores  
Dr Marc Van der Kamp  
Bioinformatics 2023

## 1   **Abstract**

2   This study evaluates structural similarity-based search methods for identifying Diels-Alder  
3   (DA) enzyme homologs, focusing on spirotetronate cyclases. Traditional searches use  
4   sequence similarity, but structural methods like Foldseek and AlphaFold clusters offer  
5   potential advantages. We compared two approaches: Method A, which applies Foldseek  
6   with TM-align and 3Di alignment methods, and Method B, using AlphaFold clusters.

7   Our analysis found that Foldseek's performance is influenced by the choice of alignment  
8   method, with TM-align and 3Di yielding different results. Although unique hits were often  
9   irrelevant, both methods successfully identified relevant DA homologs. Method B, using  
10   AlphaFold clusters, proved more reliable, consistently retrieving accurate hits due to its  
11   structured clustering approach.

12   The study also revealed limitations in relying solely on structural similarity. Incorporating  
13   sequence similarity improved search comprehensiveness, as sequence-based methods  
14   identified additional relevant hits not captured by structural searches. Predicting closed  
15   conformations for DA enzymes using AlphaFold was challenging, highlighting biases in  
16   training data and difficulties with prediction parameters.

17   In summary, combining structural and sequence-based methods enhances the discovery of  
18   DA enzyme homologs. Effective protocols should integrate Foldseek with a well-chosen  
19   alignment method, incorporate sequence similarity for broader coverage, and use AlphaFold  
20   clusters to refine results. Future research should further explore these methods, including  
21   additional parameters and experimental validations to confirm potential DA enzyme  
22   candidates.

23

24

25

26

27

28

29

30

31

32

33

34

35

36

## 37 **Acknowledgements**

38 I would like to express my sincere gratitude to Dr. Marc Van der Kamp for his exceptional  
39 guidance and support throughout this project. His patience, insightful advice, and  
40 encouragement have been invaluable in helping me learn and discover my passion for this  
41 research. To my amazing friends: Isabella, David, Camila, Danyra and all my housemates.  
42 To Calum for being a great listener and being patient. Thank you for being my home and  
43 supporting me in this adventure. To my parents and my sister for being my main inspiration  
44 and support. To my mom, I do not have words to express my gratitude to you. Finally, to  
45 Remigio my main emotional supporter.

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

### 73    **Author's Declaration**

74    I declare that this thesis has been composed solely by myself and that it has not been  
75    submitted, in whole or in part, in any previous application for a degree. Except where states  
76    otherwise by reference or acknowledgment, the work presented is entirely my own.

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102	<b>ABSTRACT.....</b>	1
103	<b>ACKNOWLEDGEMENTS .....</b>	2
104	<b>AUTHOR'S DECLARATION.....</b>	3
105	<b>1. INTRODUCTION .....</b>	6
106	<b>2. METHODS.....</b>	10
107	<b>2.1 DATA COLLECTION .....</b>	10
108	<b>2.1.2 SEARCHING BASED ON STRUCTURE.....</b>	10
109	<b>2.1.2 SEARCHING BASED ON SEQUENCE.....</b>	10
110	<b>2.2 VISUALIZING AND SELECTING RELEVANT HITS.....</b>	10
111	<b>2.3 EVALUATING THE RELEVANT HITS 3D STRUCTURE.....</b>	11
112	<b>2.4 EVALUATING THE ACTIVE CAVITY SITE.....</b>	11
113	<b>2.4.1 ASSESSING THE PRESENCE OF THE SALT BRIDGE .....</b>	11
114	<b>2.4.2 PREDICTING THE CLOSED CONFORMATIONS.....</b>	11
115	<b>2.7 CALCULATING THE CAVITY VOLUME.....</b>	11
116	<b>3. RESULTS.....</b>	12
117	<b>3.1 STRUCTURE-BASED SEARCH RESULTS.....</b>	12
118	<b>3.1.1 METHOD A.....</b>	12
119	<b>3.1.2.1 EFFECT OF THE ALIGNMENT .....</b>	12
120	<b>3.1.2.2 EFFECT OF THE QUERY.....</b>	12
121	<b>3.1.2 METHOD B.....</b>	13
122	<b>3.1.2.1 EFFECT OF THE QUERY .....</b>	13
123	<b>3.1.3 METHOD C&amp;D.....</b>	13
124	<b>3.2 RELEVANT HITS SELECTION AND RETRIEVAL.....</b>	14
125	<b>3.3 STRUCTURE EVALUATION.....</b>	14
126	<b>3.4 ACTIVE CAVITY SIZE EVALUATION.....</b>	15
127	<b>3.4.1 SALT BRIDGE EVALUATION.....</b>	15
128	<b>3.4.2 CLOSED CONFORMATION PREDICTION.....</b>	15
129	<b>3.4.3 CAVITY VOLUME EVALUATION.....</b>	16
130	<b>4. DISCUSSION.....</b>	18
131	<b>5. BIBLIOGRAPHY .....</b>	24
132	<b>6. TABLES AND FIGURES.....</b>	31
133	<b>6.1 TABLES.....</b>	31
134	<b>6.2 FIGURES .....</b>	44
135	<b>7. APPENDIX .....</b>	49
136	<b>7.1 SUPPLEMENTARY TABLES.....</b>	49
137	<b>7.2 SUPPLEMENTARY FIGURES.....</b>	55

138

139

140

141

142

143

144

145

146

147

## 148 **Contents of Tables and Figures**

### 149 *Tables*

<b>Table Number</b>	<b>Title</b>	<b>Page</b>
1	Detailed Search Protocol: Queries, Methods, and Databases.	25
2	Analysis of Sequence Similarity Among Queries: Insights and Comparisons.	26
3	Relationship Between Sequence Similarity and Unique Hit Retrieval.	27
4	Unique Hit Distribution Across Queries: Foldseek Method A Analysis.	28
5	Query Distribution Across Clusters: Method B and Search-Derived Additions.	29
6	Unique Hit Retrieval and Positive Control Detection: Method A vs. Method B.	30
7	Comparative Effectiveness: Structure vs. Sequence in Hit Retrieval.	31
8	Comparative Analysis of Salt Bridges at AbyU-Equivalent 78-122 Positions.	32-34
9	Evaluation of Hit Conformations and Their Representation in the PDB.	35-37
10	Template-Based Closed Conformation Analysis: AbyU and Cyc15 in ColabFold.	38
11	Optimized ColabFold Predictions: Assessing Structural Hit Conformations.	39-40
12	Comparative Cavity Volume Assessment of Closed Conformation Hits Using CastP.	41

150

### 151 *Figures*

<b>Figure Number</b>	<b>Title</b>	<b>Page</b>
1	Overview of unique hits and cluster representatives retrieved by alignment methods.	42-43
2	Comparison of Sequence Similarity Networks (SSNs) for Method A and Method B.	44
3	Structural visualization of AbyU and Cyc15 with emphasis on cavities and salt bridges..	45
4	Predictions of closed conformations for A0A1Y0RPR1 using templates AbyU and Cyc15 with default parameters.	46
5	Successful closed conformation prediction of A0A345XXX2 using AbyU as a template	47
6	Prediction of the closed conformation of A0A1Y0RPR1 through variations in MSA depth and number of seeds.	48

152

153

154

155

156

## 1. Introduction

In the pharmaceutical industry, carbon-carbon bond formation is essential for synthesizing complex molecules with stereospecificity (1). The Diels-Alder (DA) reaction, a [4+2] cycloaddition between a diene and a dienophile, facilitates this by forming two carbon-carbon bonds and generating cyclohexene structures. Despite its utility, the DA reaction faces challenges such as low efficiency at room temperature and difficulties in controlling stereoselectivity, often requiring high temperatures or pressures that are energetically costly (2,3).

165

To enhance the efficiency and specificity of the DA reaction, researchers are exploring biocatalysts like enzymes. Enzymatic reactions typically exhibit higher reaction rates and specificity while operating under milder conditions. Among these, the spirotetrone cyclase family has shown promise in catalyzing true [4+2] cycloadditions (4).

170

Applying enzymes to perform the DA reaction in industry can greatly improve its performance, both in efficiency and specificity. The search for naturally occurring DA enzymes has a long history. However, only a few types of enzymes have been proven capable of catalyzing true [4+2] cycloadditions. One of the proven enzymes is the spirotetrone cyclase family. They are a group of DA enzymes involved in the synthesis of polyketide natural products with a spirotetrone moiety. Well-characterized members of this enzyme family include Pyrl4, AbyU, AbmU and AbnU (5–8).

178

For instance, AbyU is part of the abyssomicin C biosynthetic pathway. It consists of two eight-stranded antiparallel  $\beta$ -barrels with (+1)8 topology sealed on both ends. A salt bridge formed by Glu (78) and Arg (122) seals one end. At the opposite end, there is a hydrophobic loop. It is formed by a  $\beta$ 1- $\beta$ 2 linker capping loop. The central part of the barrel forms a large hydrophobic cavity that can be accessed through the capping loop. The capping loop is flexible, allowing substrate entry into the cavity. It becomes more ordered once the substrate is bound. Therefore, AbyU can adopt two conformations: open and closed. The conformation is determined by the position of the capping loop (6,9).

187

Contrary to other enzymes, there is no evidence of highly conserved amino acid features in the active site. The reaction likely proceeds through hydrophobic and steric complementarity, driven primarily by the accommodation of the substrate inside the cavity (2). To date, over 400 compounds have been identified as products of the DA reaction, but only a few DA

192 enzymes have been fully characterized. This suggests that many DA enzymes remain to be  
193 discovered.

194

195 Given the industrial importance of Diels–Ader (DA) enzymes, there is a growing interest in  
196 discovering more of these enzymes. Protein databases contain many protein sequences, but  
197 few structures have been characterized. To address this, sequence similarity searches are  
198 commonly used. They identify homologous sequences to a desired query. However, inferring  
199 functional similarity from sequence similarity is not straightforward (10). This process  
200 becomes even more complex when dealing with remote homologs. These are proteins that  
201 share structure and function but exhibit low sequence identity (11).

202

203 This complexity is evident in the case of the spirotetronate DA enzymes, which often share  
204 less than 25% sequence identity (6). Despite this, the spirotetronate cyclase Cyc15 was  
205 successfully identified using sequence-based searching. The strategy involved creating a  
206 library of sequences encoding known or predicted DA enzymes. Then, the search was  
207 expanded using BLASTp to find homologous sequences in public databases. To visualize  
208 the results, a sequence similarity network (SSN) was constructed. A SSN groups sequences  
209 into nodes connected by edges based on a selected sequence similarity threshold. A  
210 representative sequence from each node was then expressed and tested for enzymatic  
211 activity, using the natural substrate of AbyU. Notably, Cyc15 exhibited considerable DA  
212 activity and resulted in a product with the opposite stereochemistry to the AbyU product (6).  
213 This case demonstrates the effectiveness of using sequence similarity to find homologous  
214 proteins. However, it may not be comprehensive: some hits can be missed. In the case of  
215 DA enzymes, there is no evidence of specific sequence features retained in the active site.  
216 Consequently, local and global sequence identity may be too small to be statistically  
217 significant (12).

218

219 Structural alignment tools (structure-based methods) offer an alternative approach for  
220 detecting homologous proteins with similar functions. Such methods compare the overall fold  
221 of the query protein to identify hits that share the same structure. Another alternative is  
222 protein language models. They represent proteins as embeddings. They can capture  
223 complex patterns in sequences, such as PLM-Blast (13).

224

225 Foldseek is an efficient structural search algorithm that converts the query to the dictionary  
226 which describes the tertiary interactions of amino acids (3Di dictionary). This allows for the  
227 identification of proteins with a common structure (14). In the Foldseek web server, a 3Di-  
228 based search in structural databases can be performed, with the hits ranked based on the e-

229 value obtained. As an alternative to 3Di, TM-align can be used. It performs the alignment  
230 based on the position of the backbone C $\alpha$  atoms (15). Foldseek has been used to cluster the  
231 AlphaFold database (AFDB), creating a database called AlphaFold Clusters (16). The  
232 COFACTOR platform offers another approach to searching protein structures based on the  
233 TM-align algorithm, using the BioLip database (17)

234

235 These methods provide alternatives to sequence-based searching for finding distant  
236 homologs. Thus, they can be applied to expand the search for spirotetronate DA enzymes.  
237 Foldseek has demonstrated high efficiency, sensitivity, and alignment quality. A benchmark  
238 study compared its performance with other structural aligners such as Dali, CE, and TM-  
239 align (14). It concluded that Foldseek sensitivity is comparable to Dali and significantly  
240 higher than CLE-SW and MMseqs2. These findings support the use of Foldseek for finding  
241 spirotetronate DA homologs.

242

243 Many proteins share the  $\beta$ -barrel fold found in the characterized spirotetronate DA enzymes.  
244 However, they do not all display DA (18). Therefore, additional features must be evaluated to  
245 determine whether a hit could be a DA homolog. One key feature is the presence of a salt  
246 bridge. Such an evaluation should include the presence of a suitable cavity for the substrate  
247 to be confined in.

248

249 To accurately compare cavity metrics, DA enzymes must be in their closed conformations.  
250 This requires the capping loop to adopt its ordered structure and seal one end of the  $\beta$ -  
251 barrel. AlphaFold2 (19), e.g. as applied through ColabFold (20), is a tool capable of accurate  
252 protein conformation predictions and could thus be used to predict the closed conformation  
253 of DA enzymes.

254

255 Structure-based approaches offer a promising strategy for detecting remote homologs.  
256 However, their implementation could have some challenges. First, there is no previous  
257 structure-based search protocol successfully tested for finding DA enzymes. Therefore, the  
258 main aim of this study is to develop an effective and reliable search protocol to identify  
259 putative spirotetronate DA enzymes. The different parameters (query, alignment method) will  
260 be assessed. Previously characterized members of the spirotetronate DAse family will be  
261 used as queries. The results obtained by the different methods (with their specific  
262 parameters) will be compared. A Sequence Similarity Network (SSN) will be used to  
263 differentiate hits that could be relevant. Another important aim of this study is to find DA  
264 enzymes with alternative properties (e.g. different cavity size or stereoselectivity). Therefore,  
265 the relevant hits obtained will be structurally evaluated to compare them with current known

266 DA enzymes. Finally, predictions of their closed conformation will be performed to evaluate  
267 their cavity volume.

268

## 269      2. Methods

270      Note: all methods were run in default parameters unless otherwise indicated.

271      2.1 *Data collection*

272          2.1.1 *Searching based on structure*

273      Different structural-based tools were used to find DA enzyme homologs. They were  
274      chosen according to previous benchmark studies that proved their effectiveness and  
275      sensitivity for finding remote homologs (14,16). Foldseek ([search.foldseek](#)) was used for  
276      method A using TM-align and 3Di search. AlphaFold clusters ([search.AFDBcluster](#)) were  
277      used for method B. The clusters obtained by each query were collected. The web server  
278      also provides a list of 10 similar clusters every search. The ten similar clusters to AbyU  
279      were included to complement the results. Structural alignment based on Ca positions  
280      has also proven high sensitivity, although at the expense of speed. To have a robust  
281      comparison between current structural-based search tools, COFACTOR  
282      ([search.COFACTOR](#)) was used for Method C. Finally, protein language models are a  
283      new alternative to structural alignments. PLM-blast has been reported to excel at distant  
284      homolog proteins search (21). Therefore, Method D employed PLM-blast ([search.PLMBlast](#)). All search methodologies are detailed in Table 1. Hits sharing >95% identity  
285      (based on Blastp) were discarded (Table S1). The presence of previously characterized  
286      spirotетronate cyclases with confirmed or predicted DA activity(6) was evaluated  
287      (positive control) (Table S2). PDB files and fasta sequences of each query used can be  
288      found in Figure S1 and Table S3.

291          2.1.2 *Searching based on sequence*

292      The fasta sequences obtained by (6)were collected and searched in UniProt  
293      ([search.UniProt](#)). The variants (sequences with >95% identity) were discarded using  
294      Blastp. The sequences that did not yield matches in UniProt were discarded. These hits  
295      are named ‘sequence hits’.

297      2.2 *Visualizing and selecting relevant hits*

298      Sequence similarity networks (SSNs) were built using EFI-EST (22) with UniProt IDs.  
299      from each search method (sequences without start or stop codons were excluded).  
300      Table S4 lists the parameters used for each network. The SSNs were visualized with  
301      Cytoscape (23), using the organic layout. Only members of the cluster(s) containing the  
302      positive controls were selected for further analysis (name as relevant hits).

304      2.3 *Evaluating the relevant hits 3D structure*

305 Structure visualization was performed in PyMOL (24), using structures from PDB if  
306 present, or else from the Alphafold database ([search.AlphaFold](#)). Alignment to AbyU was  
307 performed in PyMOL. For sequence hits do not present in UniProt, structures were  
308 predicted using ESMfold ([searchESMMetagenomicAtlas](#)) (25).

309

310 *2.4 Evaluating the active cavity site*

311 *2.4.1 Assessing the presence of the salt bridge*

312 Hits were aligned based on sequence using Clustal Omega (26) ([search.ClustalOmega](#)).  
313 The MSA was visualized using a guide phylogenetic tree made by Clustal Omega. They  
314 were aligned based on structure using TExpresso (27) ([searchTCoffee.com](#)).

315

316 *2.4.2 Predicting the closed conformations*

317 Predicted structures were considered to have a closed conformation if the capping loop  
318 acquired a similar position to that of AbyU (see Figure S1). AF2 prediction in ColabFold  
319 (20) was used in an attempt to predict the closed conformation of remaining hits.

320 Different strategies were tested:

- 321 - Closed structures of AbyU ([PDB: 5DYQ](#)) or Cyc15 (AFDB prediction of [UniProt](#)  
322 [A0A401MXE6](#)) were used as a template (template-based). The highest-scoring models  
323 were visualized.
- 324 - Adjustments were made to dropout rates, number of seeds (num\_seeds), and MSA  
325 depth (max\_msa) (parameter-based). The first four highest-scoring models were  
326 visualized.

327

328 *2.4.3 Calculating the cavity volume*

329 CastP ([searchCastp](#)) (28) was used to calculate the cavity volume using a sphere radius  
330 185 size of 1.1Å.

331

### 332      3. Results

#### 333      3.1 Structure-based search results

##### 334          3.1.1 Method A

###### 335            3.1.2.1 Effect of the alignment

336      Method A employed Foldseek with four different queries and two alignment methods (3di or  
337      TM) (see Table 1). The hits obtained are summarized in Table 2. In general, both alignment  
338      methods retrieved a similar number of hits independent of the query. However, each  
339      alignment method retrieved some unique hits that were not found by the other. If these  
340      unique hits are irrelevant (e.g., not complete  $\beta$ -barrels structures) then the search can be  
341      performed either using TM-align or 3Di. To confirm this, the unique hits were searched in the  
342      list of relevant hits (sub-cluster members obtained after SSN visualization). None of them  
343      were found, thus the relevant hits were found by both methods equally. The highest-scoring  
344      hits retrieved by TM-align and 3di were evaluated (Figure 1A). They differ significantly in  
345      structure. B6TGG4 (retrieved by TM-align) contains a partial  $\beta$ -barrels with a loop, while  
346      Q75KD7 (retrieved by 3Di) contains a  $\beta$ -barrels with an additional non-folded structure.

347

###### 348            3.1.2.2 Effect of the query

349      The number of (unique) results also differs per query, with Cyc15 recovering the most  
350      unique hits. A Blastp alignment was conducted to assess the similarity between the queries  
351      (all vs. all). The goal was to determine if this similarity affects the number of unique hits  
352      retrieved. The logic was that the more similar the queries are to each other, the fewer unique  
353      results they should retrieve. AbnU shares 49% similarity with AbyU (see Table 3) when  
354      compared directly, yet each still retrieved unique results. To further assess this, a similar  
355      sequence to AbyU (identified in the phylogenetic tree in Figure S3) was used as a query in  
356      Foldseek. The query chosen, A0A246RFQ9, shares 83% similarity with AbyU. The hits  
357      obtained were compared to those from AbyU. A0A246RFQ9 retrieved some unique hits, but  
358      to a much lesser extent compared to AbnU (Table 4). This suggests that considering  
359      sequence similarity between queries is relevant to expanding the search. However, the  
360      unique hits retrieved by each query were also searched in the list of relevant hits and they  
361      were ultimately irrelevant for this study. Then, using different queries expands the search but  
362      does not increase its specificity as all the relevant hits were found by all the queries. The  
363      unique hits still display structures  $\beta$ -barrel like structures, but they are incomplete, or they  
364      are part of a larger complex. The highest-scoring unique hits retrieved by each query were  
365      visualized in PyMOL (Figure 1B). Their respective score is also shown (either tm or 3di). For  
366      reference, a hit with a TM score below 0.5 or and e-value higher than 0.01 is often  
367      considered not significant (14).

368  
369 Notably, method A did not retrieve AbyU as a hit when AlphaFold/Uniprot50 v4,  
370 AlphaFold/Swiss-prot v4, AlphaFold/Proteome v4 were used as databases. The search was  
371 expanded by adding PDB100. After performing these changes, AbyU was found. A possible  
372 explanation could be that the pre-filter stage of Foldseek was too restrictive. The pre-filter  
373 performs an alignment that reduces false positives and non-homologous sequences. The  
374 criterion is based on double k-mer matches which non-homologous sequences are less  
375 likely to have (14). However, when turning off the pre-filter AbyU was not retrieved. Other  
376 reasons not explored here influenced Foldseek's inability to retrieve AbyU when the PDB100  
377 is not used as a database.

378

### 379        3.1.2 Method B

380        3.1.2.1 Effect of the query  
381 The hits obtained by method B are summarized in Table 5. AbnU, AbyU and Cyc15 were  
382 clustered together, the total members of the cluster were 78 proteins. Pyrl4 was clustered  
383 separately with 32 proteins. The total results obtained are higher compared to the hits  
384 retrieved by method A. However, all the hits are concentrated in the A0A022Q700 cluster  
385 with 4000 members. Notably, the effect of the query is very different for both methods. In  
386 method A each query retrieved a significant amount of unique hits. Conversely, in method B  
387 AbnU, AbyU and Cyc15 retrieved the same hits. Cluster A0A0K9XAS6 and A0A3A9ZPZ1  
388 were found by the queries. Therefore, the members of these clusters could be relevant hits.  
389 They were searched in the list of Method B's relevant hits (121 hits). Clusters A0A0K9XAS6  
390 and A0A3A9ZPZ1 retrieved 75 hits of the relevant hits. The other 34 hits were found by the  
391 added clusters. In that sense, including the additional clusters expands the search and still  
392 retrieves relevant hits. The search is enriched when queries with lower sequence similarity  
393 (Pyrl4) are added. The clusters' representative structure was evaluated in Pymol (Figure  
394 1C). They all display a  $\beta$ -barrels structure that resembles the queries used. Therefore,  
395 expanding the search in method B proved to be more effective in this study than in method  
396 A.

397

### 398        3.1.3 Method C & D

399 The hits obtained method C and D were also evaluated (Table S6 and Table S7). Currently  
400 these web servers only employ PDB100 as the database. Therefore, the number of hits  
401 obtained are significantly lower compared to the other methods. The hits were evaluated  
402 manually. Method C was originally run as a "state-of-the-art" approach, where all the positive  
403 controls were expected to be found. As expected, method C hits included well-known DA  
404 enzymes. However, it retrieved some other hits that are not characterized by DA enzymes,

405 such as pterocarpan synthases (29). Also, an allene oxide cyclase is involved in oxylipin  
406 biosynthesis (30). The hits obtained by method D included some well-known spirotetronate  
407 DA enzymes. However, other hits displayed non-beta barrel structures and beta barrels  
408 without DA activity. For instance, JAC1 ([JAC1.PDB](#)) has been reported to have chaperone  
409 activity (31). Therefore, the results obtained by methods C and D were not considered for  
410 further analysis, except for those already included in other methods.

411

### 412 *3.2 Relevant hits selection and visualization*

413 After the results retrieval for each method, the SSNs were visualized. The SSN of methods A  
414 and B were compared. The connection threshold used was 25% sequence similarity. For  
415 both methods, all the positive controls were clustered together and in a separated cluster.  
416 Therefore, members of these sub-clusters were considered relevant hits. The relevant hits  
417 unique to each method were counted (Table 6). The sub-clusters of each method can be  
418 found in Figure 2. Both methods retrieved members of the positive control. Method B,  
419 retrieved the most compared to method A. However, is important to consider various search  
420 methodologies as both methods found unique hits.

421

### 422 *3.3 Structure evaluation*

423 Given that they were clustered together, all the relevant hits were expected to display an 8-  
424 stranded β-barrels. However, some of the structures found were not β-barrels, they  
425 displayed different folds. K4FDH3 displayed a structure containing two beta barrels  
426 connected by a loop with different sequence content (Figure S4A)., Each beta barrel was  
427 thus evaluated as an individual hit. A0A522WLH4 and A0A838IMS4 displayed a structure  
428 indicated as ‘complex’. It contained a β-barrels similar to other structures, but also significant  
429 additional structural elements. Thus, they were not included in further analysis (Figure S4B  
430 and Figure SBC).

431

432 The structures of the hits previously obtained using sequence-based search (6)were also  
433 evaluated. Some of the hits displayed a β-barrels structure but other hits displayed other  
434 folds. The combined results from the structure-based search (methods A and B) were  
435 counted and compared with the hits obtained by Zorn et al (see Table 7). The structure-  
436 based search retrieved more unique hits, although both methods obtained hits that were not  
437 found by the other. Therefore, searching homologs based on sequence similarity also  
438 retrieves relevant hits. However, searching by structure provided more results that displayed  
439 the β-barrels fold. Only hits obtained by structure-based search (structural hits) were  
440 considered for further analysis.

441

442     *3.4 Active Cavity size Evaluation*443         *3.4.1 Salt bridge evaluation*

444     Most of the reactions in the pharmaceutical industry are inter-molecular. The ideal DA  
445     enzyme will have enough cavity space to perform a reaction with two substrates. AbyU  
446     typically displays a salt bridge formed by Arg-His in positions 78 and 122 (see Figure 3A)  
447     The presence of the salt bridge limits the cavity volume, thus hits without the Arg -His salt  
448     bridge will be of interest. Structure-based sequence alignment (using T-COFFEE) was  
449     performed. The presence of a salt bridge that caps the active site (positions 78 and 122 in  
450     AbyU numbering) was found in many, but not all hits (Table 8). Some hits do not contain the  
451     same Glu/Asp-Arg amino acid pair (as in AbyU and others) but have similar amino acids.  
452     They can still cap the cavity through hydrogen bonding (e.g. Arg replaced by Tyr) Further,  
453     bulky amino acids like Arg would still affect the cavity volume even if they do not form a salt  
454     bridge with the amino acid on the opposite side of the barrel. Therefore, to reveal hits with  
455     possible alternative cavity shapes/sizes, particular attention was given to hits that do not  
456     contain Arg. However, despite not having Arg in the same position, some hits still have bulky  
457     amino acids such as leucine (L), methionine (M), tyrosine (Y), tryptophan (W), and leucine  
458     (L). Their presence could still limit the cavity volume even without forming a salt bridge. In  
459     position 122, the amino acids tend to be bulky. Conversely, in position 78 the amino acids  
460     vary significantly in composition and do not follow a particular trend.

461

462         *3.4.2 Closed conformation prediction*

463     As described in the introduction for AbyU, the cavity is often delimited by a salt bridge at one  
464     end and the capping loop at the other. Then, the cavity is lined by the side chains of several  
465     residues that point into the cavity (Figure 3). When the loop acquires an ordered structure,  
466     the proteins are considered in its closed conformation. The conformations present in the  
467     AFDB either have an open or a closed capping loop or domain (Table 9).

468     In an attempt to obtain closed conformations for those hits with an open conformation,  
469     custom predictions were tested. Tests were initially performed with hit A0A1Y0RPR1  
470     ([AFDBentry](#)), using closed structures of AbyU (from PDB ID 5DYV) or Cyc15 (AlphaFold  
471     prediction of A0A401MXE6) as templates (Figure 4 and Table 10). When a custom template  
472     is used, this structure is used as the basis for the initial ‘pair representation’ prediction in  
473     AlphaFold. It was assumed that if the template and the query share a higher sequence  
474     identity in the capping loop region, it will help get alternative conformations. However, with  
475     both AbyU and Cyc15 as templates (with Cyc15 being more similar to A0A1Y0RPR1 in the  
476     loop region), an open conformation was predicted (Figure S6). Notably, only one temple-  
477     based prediction was successful, A0A345XXX2 ([AFDB.entry](#)) using AbyU as a template  
478     (Figure 5).

479 Therefore, further prediction attempts used a different approach. They were based on  
480 previous reports that successfully predicted alternative protein conformations. The  
481 assumption was that conformations can be decoded from sequence (32)Using the Da Silva  
482 approach, the MSA and the number of seeds were changed. Alternative capping  
483 conformations were predicted (Figure 6, Table 11), but this does not resemble a fully closed  
484 conformation.

485

486 Notably, the alternative conformations were obtained when the number of seeds was  
487 changed to 16 in all the MSA depths. Then, it was hypothesized that the number of seeds  
488 was controlling the change in the conformations rather than the MSA depth.

489 Additional predictions were made while keeping the MSA depth parameter in its default  
490 setting and varying the number of seeds. All the number of seeds (from 2-8) obtained some  
491 movement of the capping loop (except for number of seeds of 16). When compared to the  
492 previous predictions where the MSA depth was also changed (Figure 11C), the movement of  
493 the capping loop is less evident. Therefore, changing the MSA depth and the number of  
494 seeds had a stronger effect on the capping loop position.

495

#### 496       3.4.3 Cavity volume evaluation

497 Although closed models were not obtained for all hits, the cavity of those models that do  
498 have a closed conformation can be further evaluated. Hits without the salt bridge are  
499 expected to have a higher cavity volume, as the cavity may extend further into the barrel (as  
500 observed for Pyrl4). After separating hits that either have an Arg at position 122 ('salt bridge'  
501 or similar capping of the cavity) or not, mean cavity volumes of the closed hits were  
502 calculated (Table 12). The results indicate that in general, hits lacking the salt bridge tend to  
503 have a higher cavity volume. However, to attribute this to the presence of the salt bridge is  
504 difficult to conclude. Only 4 cavity volumes without the salt bridge were measured (compared  
505 to 17 with the salt bridge). Therefore, there is not enough data to make a confident  
506 comparison.

## 507 4. Discussion

508 This study investigated different search protocols based on structural similarity to find DA  
509 homologs. While earlier studies have successfully found DA enzymes using sequence  
510 similarity searches, they have not employed structural similarity as a search parameter.  
511 Given the low sequence similarity they share, using structural-based search may help  
512 uncover new DA enzymes. Foldseek (method A) and Alpha Fold clusters ( method B) were  
513 employed as search algorithms. As there are no previous reports using these methods to  
514 find DA enzymes, the parameters of each search were compared. The goal is to develop an  
515 effective search method that expands the search but also is specific enough to find putative  
516 DA enzymes. Additionally, the relevant hits obtained were evaluated structurally (cavity  
517 measurement) to present them as a possible DA enzyme which could later be tested  
518 experimentally.

519

520 We found that the Foldseek search is more complex than expected. Initially, Foldseek was  
521 thought to retrieve the same hits independent of the query or alignment method. In Method A  
522 each query and alignment method retrieved unique hits (see Table 2). However, the unique  
523 hits were ultimately irrelevant. In a benchmark study (14), different alignment methods were  
524 compared including Foldseek (3Di) and Foldseek-TM. Their ability to classify members of  
525 the same family and superfamily was tested (Scope dataset). Both methods had the highest  
526 and third highest performance. Although very similar, they still had small variations when  
527 finding true positives. Therefore, when used in the webserver they are expected to retrieve  
528 different results. One of the key differences between the alignment methods relies on the  
529 local and global differences detection. TM-align specializes in detecting global similarities. In  
530 contrast, 3Di focuses on local structures. However, both methods utilize the 3Di dictionary to  
531 convert the query. TM-align, then, is not applied in its 'pure' form as in other algorithms (eg.  
532 COFACTOR). The 3Di conversion considers tertiary interaction between aminoacidic pairs,  
533 while TM-align only considers C $\alpha$  positions. The different hits retrieved by each method,  
534 then, could be due to the different scoring systems that each alignment applies. However,  
535 the effect of the scoring system on the structures of the hits retrieved is not evident. Overall,  
536 each alignment method retrieved unique hits that do not have an evident structural  
537 difference between them. The alignment method, then, does not capture a specific structure.  
538 For this study, both alignment methods effectively retrieved the relevant hits. Therefore,  
539 when using method A the search can be performed with either TM-align or 3Di indistinctly.

540

541 The queries used in Foldseek also impact the hits obtained. Given that they share a similar  
542 fold, they were expected to retrieve similar hits. However, each query retrieved different

543 results. This divergence could be explained by how Foldseek converts aminoacidic  
544 composition into the 3di dictionary. For each residue i, it identifies its nearest neighbour j  
545 which is then categorized into one of 20 possible states. Although the descriptors consider  
546 spatial conformation, the amino acid identities also influence how these conformations are  
547 described. Therefore, even though the queries AbyU, AbnU, Cyc15, and PyrL4 share a  
548 similar fold, their sequence content influences the retrieval of different results. This is shown  
549 in Table S4, A0A246RFQ9 still manages to retrieve unique hits but at a much lower capacity  
550 than AbnU. However, the unique hits were irrelevant. When comparing their structure, they  
551 displayed a variety of folds. Notably, their respective scores (TM or e-value) are much lower  
552 when compared to the relevant hits.

553 All the queries retrieved the relevant hits, irrespective of their sequence content. Conversely,  
554 non-relevant hits display different folds. This suggests that the highest information density is  
555 encoded in conserved protein cores, while non-conserved coil/loop regions carry lower  
556 information density (14). The 3Di effectively captures this, as all the highest-scoring hits  
557 were found by all queries and most of them displayed the  $\beta$ -barrels. Conversely, the unique  
558 hits retrieved by each query were retrieved when aligning partial  $\beta$ -barrels structures or  
559 additional motifs which were found by the specific sequence content of each query.  
560

561 Moreover, we found that relying uniquely on structure is not effective. The importance of  
562 sequence similarity was first shown in the different hits obtained by methods A and B. Both  
563 methods retrieved members of the positive control, but method B proved to be more  
564 effective. It retrieved more unique hits and positive control than method A (see Table 4).  
565 Another key difference was the impact of the query in both methods which was not  
566 expected. Although Foldseek and AlphaFold clusters are different algorithms, they both  
567 employ the 3Di dictionary. Then, their implementation should be similar. In the AFDB AbyU,  
568 AbnU, and Cyc15 are clustered together. When used as queries in foldseek, then, they  
569 should retrieve similar hits. However, each query retrieved a large amount of unique hits.  
570 This is partly explained by the difference in sequence content. However, it does not fully  
571 explain why in method A queries (AbyU, AbnU, and Cyc15) could expand the search while in  
572 method B they retrieved the same results. This divergence can be explained by the  
573 difference between the methods. Unlike the Foldseek algorithm, the initial stage of clustering  
574 for the AF clusters is based on a multiple sequence alignment (MSA). It uses MMseqs with  
575 50% identity and 90% overlap to select a cluster representative. Subsequently, the 3Di  
576 dictionary from Foldseek is used for structure-based clustering (16). AF clusters filter hits by  
577 a sequence identity threshold, whereas Foldseek skips this step. As a result, in method A  
578 the query is compared with a larger pool of target sequences. Then, the hits obtained are  
579 much more diverse than method B. However, this diversity does not mean that the hits are

580 relevant. Adding more queries in method B increases the relevant hits obtained. Conversely,  
581 in method A, as there is no previous sequence identity filter adding more queries did not  
582 increase the relevant hits retrieved.

583

584 Sequence similarity also proved to be useful when filtering relevant results. In the SSN of all  
585 methods, the relevant hits formed a separated cluster. Furthermore, sequence-based search  
586 also retrieved results that were not found by structural-based search. Therefore, sequence  
587 similarity needs to be considered to perform a reliable search. Pereira et al. (33) showed the  
588 importance of a combined approach. They successfully classified the ‘dark matter’ proteins  
589 in UniProt by using structure features and sequence similarity. Therefore, a combined  
590 approach with sequence similarity and structural alignment proves to be effective for finding  
591 remote DA homologs.

592

593 Predicting the closed conformation proved more challenging than anticipated. Notably, most  
594 of the hits did not need a prediction as they were already in a closed conformation. Del  
595 Alamo et al. ((34), concluded that AlphaFold is biased towards known structures. This can  
596 be explained by the data that was used to train AlphaFold neural networks. The data  
597 available in PDB before April 30 2018 was used. By 2018, only AbyU and Pyrl4 from the  
598 positive control had a crystal structure entry in PDB. AbyU has 5DYQ ([open conformation](#))  
599 and 5DYV ([closed conformation](#)) and Pyrl4 has 5BTU ([open conformation](#)) and 7DVK  
600 ([closed conformation](#)). Thereby, there is a bias toward previously known proteins (34).  
601 AlphaFold already has the information of two characterized DA conformational states (open  
602 or closed). Then, the structure prediction of putative DA enzymes could acquire the open or  
603 closed conformation. However, using ColabFold to predict the closed conformation of the  
604 open-structured hits was challenging. The use of templates in their closed conformation  
605 (AbyU or Cyc15) was irrelevant to predicting an alternative conformation. In contrast,  
606 changing the MSA depth (max\_msa parameter) predicted intermediate conformations. This  
607 agrees with the results obtained by Del Alamo et al. where a template in the outward-facing  
608 (OF) conformation of MCT1 was provided to ColabFold, along with MSA of various sizes.  
609 The desired OF conformation was obtained only with MSAs of 16-32 sequences. The other  
610 conformations obtained were either inward-facing (IF) or intermediate conformations. The  
611 same results were observed when using LAT1 in its OF or IF template or PTH1R in its active  
612 or inactive state. Their alternative conformations were obtained only when the MSA was  
613 shallow (34). This suggests that templates are useful when the MSA depth is changed. As  
614 the MSA depth increases, the template becomes less informative. Moreover, as suggested  
615 by Wayment-Steele et al. different proteins need an appropriate MSA depth to predict a  
616 desired conformation (35). In this case, the accurate MSA depth was not found. The number

617 of seeds was also a relevant parameter to obtain different conformations. Then, the  
618 appropriate parameters combination (num\_seeds and max\_seq) would predict a closed  
619 conformation.

620

621 This study, however, has some limitations. First, when accessing the difference between the  
622 methods (alignment and queries) only the highest scoring hits were evaluated. To have a  
623 robust comparison, all the unique hits should be evaluated. Their respective scores and  
624 structures should be compared. However, as they were ultimately irrelevant, they did not  
625 affect the results obtained.

626 The closed conformations prediction also has some limitations. First, not many DA enzymes  
627 have been characterized. Then, using tools that rely on neural networks may limit the results  
628 obtained as they are trained on characterized proteins available in databases (e.g. [PDB](#)).  
629 Alternatively, another prediction tool could be used to predict the closed conformations. This  
630 was performed by Mbatha *et al.* (9), they used Modeller ([search\\_modeller](#)) to predict the  
631 closed conformation of AbmU. However, using ColabFold helped notice the existing AF bias  
632 in DA conformation prediction. Moreover, it also showed the importance of configuring the  
633 MSA depth and number of seeds.

634

635 Avenues for future research include evaluating a set of remote homolog enzyme functions.  
636 For instance, the complete search protocol (including search, SSN visualization, and  
637 structure evaluation) had relevant results for finding putative DA enzymes. However, other  
638 enzymes should be tested to conclude if combining structure-based search and sequence  
639 similarity is an optimal strategy. Cavity volume and the presence of the salt bridge are  
640 parameters that help classify DA enzymes. However, the evaluation of more parameters is  
641 important to evaluate alternative spiroketone DA enzymes. For instance, cavity size and  
642 cavity lining residues can be included as classification parameters. Moreover, docking  
643 experiments and reaction simulations should be implemented to state that a hit could be an  
644 alternative spiroketone DA enzymes.

645

646 While new structural-based tools to search remote homologs are available, an effective  
647 protocol for implementing these tools in the search for unknown DA enzymes has not been  
648 done yet. The parameters, the algorithm, and the sequence similarity used have an impact  
649 on the results obtained. To do an effective search, sequence similarity, and structural  
650 approaches must be combined. These are key recommendations when searching remote  
651 enzyme homologs:

652

653 - Foldseek: the search parameters can include one well-known query and any  
654 alignment method (TM-align or 3Di). The use of multiple queries is not necessary.  
655 - Sequence similarity: Alternatively, an MSA can be done with sequences retrieved  
656 from multiple databases. Then, only sequences sharing a % of sequence similarity will be  
657 kept. Next, they can be used as queries in Foldseek.  
658 - AlphaFold clusters: The first step should be identifying the clusters containing the  
659 well-known enzymes. Then, uncharacterized members of this cluster can be used as queries  
660 to search again. This step can be repeated until the structure does not resemble the initial  
661 query's fold.  
662  
663 These recommendations would help refine the search and increase the possibility of finding  
664 a putative remote homolog of the desired enzyme.  
665  
666  
667

## 5. Bibliography

1. Gao L, Yang J, Lei X. Enzymatic intermolecular Diels-Alder reactions in synthesis: From nature to design. Vol. 2, *Tetrahedron Chem.* 2022.
2. Byrne MJ, Lees NR, Han LC, Van Der Kamp MW, Mulholland AJ, Stach JEM, et al. The Catalytic Mechanism of a Natural Diels-Alderase Revealed in Molecular Detail. *J Am Chem Soc.* 2016;138(19).
3. Gregoritza M, Brandl FP. The Diels–Alder reaction: A powerful tool for the design of drug delivery systems and biomaterials. *European Journal of Pharmaceutics and Biopharmaceutics.* 2015 Nov 1;97:438–53.
4. Sara AA, Um-E-Farwa UEF, Saeed A, Kalesse M. Recent Applications of the Diels-Alder Reaction in the Synthesis of Natural Products (2017–2020). *Synthesis (Germany).* 2022;54(4).
5. Tian Z, Sun P, Yan Y, Wu Z, Zheng Q, Zhou S, et al. An enzymatic [4+2] cyclization cascade creates the pentacyclic core of pyrroindomycins. *Nat Chem Biol.* 2015;11(4).
6. Zorn K, Back CR, Barringer R, Chadimová V, Manzo-Ruiz M, Mbatha SZ, et al. Interrogation of an Enzyme Library Reveals the Catalytic Plasticity of Naturally Evolved [4+2] Cyclases. *ChemBioChem.* 2023;24(14).
7. Li B, Guan X, Yang S, Zou Y, Liu W, Houk KN. Mechanism of the Stereoselective Catalysis of Diels-Alderase PyrE3 Involved in Pyrroindomycin Biosynthesis. *J Am Chem Soc.* 2022;144(11).
8. Kashyap R, Yerra NV, Oja J, Bala S, Potuganti GR, Thota JR, et al. Exo-selective intermolecular Diels–Alder reaction by Pyrl4 and AbnU on non-natural substrates. *Commun Chem.* 2021;4(1).
9. Mbatha SZ, Back CR, Devine AJ, Mulliner HM, Johns ST, Lewin H, et al. Antibiotic origami: selective formation of spiroketones in abyssomicin biosynthesis. *Chem Sci.* 2024;
10. Pearson WR. An introduction to sequence similarity (“homology”) searching. *Curr Protoc Bioinformatics.* 2013;(SUPPL.42).
11. Chen J, Guo M, Wang X, Liu B. A comprehensive review and comparison of different computational methods for protein remote homology detection. *Brief Bioinform.* 2018;19(2).
12. Pertsemidis A, Fondon JW. Having a BLAST with bioinformatics (and avoiding BLAST phemy). Vol. 2, *Genome Biology.* 2001.
13. Kaminski K, Ludwiczak J, Pawlicki K, Alva V, Dunin-Horkawicz S. pLM-BLAST: distant homology detection based on direct comparison of sequence representations from protein language models. *Bioinformatics.* 2023;39(10).

- 704 14. van Kempen M, Kim SS, Tumescheit C, Mirdita M, Lee J, Gilchrist CLM, et al. Fast  
705 and accurate protein structure search with Foldseek. *Nat Biotechnol.* 2024;42(2).
- 706 15. Zhang Y, Skolnick J. TM-align: A protein structure alignment algorithm based on the  
707 TM-score. *Nucleic Acids Res.* 2005;33(7).
- 708 16. Barrio-Hernandez I, Yeo J, Jänes J, Mirdita M, Gilchrist CLM, Wein T, et al. Clustering  
709 predicted structures at the scale of the known protein universe. *Nature.*  
710 2023;622(7983).
- 711 17. Zhang C, Freddolino PL, Zhang Y. COFACTOR: Improved protein function prediction  
712 by combining structure, sequence and protein-protein interaction information. *Nucleic  
713 Acids Res.* 2017;45(W1).
- 714 18. Fairman JW, Noinaj N, Buchanan SK. The structural biology of β-barrel membrane  
715 proteins: A summary of recent reports. Vol. 21, *Current Opinion in Structural Biology.*  
716 2011.
- 717 19. Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O, et al. Highly  
718 accurate protein structure prediction with AlphaFold. *Nature.* 2021;596(7873).
- 719 20. Mirdita M, Schütze K, Moriwaki Y, Heo L, Ovchinnikov S, Steinegger M. ColabFold:  
720 making protein folding accessible to all. *Nat Methods.* 2022;19(6).
- 721 21. Pantolini L, Studer G, Pereira J, Durairaj J, Tauriello G, Schwede T. Embedding-  
722 based alignment: combining protein language models with dynamic programming  
723 alignment to detect structural similarities in the twilight-zone. *Bioinformatics.*  
724 2024;40(1).
- 725 22. Zallot R, Oberg N, Gerlt JA. The EFI Web Resource for Genomic Enzymology Tools:  
726 Leveraging Protein, Genome, and Metagenome Databases to Discover Novel  
727 Enzymes and Metabolic Pathways. *Biochemistry.* 2019;58(41).
- 728 23. Otasek D, Morris JH, Bouças J, Pico AR, Demchak B. Cytoscape Automation:  
729 Empowering workflow-based network analysis. *Genome Biol.* 2019;20(1).
- 730 24. DeLano WL. Pymol: An open-source molecular graphics tool. CCP4 Newsletter on  
731 protein crystallography. 2002;40.
- 732 25. Lin Z, Akin H, Rao R, Hie B, Zhu Z, Lu W, et al. Evolutionary-scale prediction of  
733 atomic-level protein structure with a language model [Internet]. Vol. 379, *Science.*  
734 2023. Available from: <https://esmatlas.com>
- 735 26. Sievers F, Higgins DG. Clustal Omega. *Curr Protoc Bioinformatics.* 2014;2014.
- 736 27. Poirot O, O'Toole E, Notredame C. Tcoffee@igs: A web server for computing,  
737 evaluating and combining multiple sequence alignments. *Nucleic Acids Res.*  
738 2003;31(13).
- 739 28. Tian W, Chen C, Lei X, Zhao J, Liang J. CASTp 3.0: Computed atlas of surface  
740 topography of proteins. *Nucleic Acids Res.* 2018;46(W1).

- 741 29. Meng Q, Moinuddin SGA, Kim SJ, Bedgar DL, Costa MA, Thomas DG, et al.  
742 Pterocarpan synthase (PTS) structures suggest a common quinone methide-  
743 stabilizing function in dirigent proteins and proteins with dirigent-like domains. *Journal  
744 of Biological Chemistry.* 2020;295(33).
- 745 30. Neumann P, Brodhun F, Sauer K, Herrfurth C, Hamberg M, Brinkmann J, et al.  
746 Crystal Structures of physcomitrella patens AOC1 and AOC2: Insights into the  
747 enzyme mechanism and differences in substrate specificity. *Plant Physiol.*  
748 2012;160(3).
- 749 31. Huwa N, Weiergräber OH, Fejzagić A V., Kirsch C, Schaffrath U, Classen T. The  
750 Crystal Structure of the Defense Conferring Rice Protein OsJAC1 Reveals a  
751 Carbohydrate Binding Site on the Dirigent-like Domain. *Biomolecules.* 2022;12(8).
- 752 32. da Silva GM, Cui JY, Dalgarno DC, Lisi GP, Rubenstein BM. Predicting relative  
753 populations of protein conformations without a physics engine using AlphaFold2.  
754 *Biophys J.* 2024;123(3).
- 755 33. Durairaj J, Waterhouse AM, Mets T, Brodiazhenko T, Abdullah M, Studer G, et al.  
756 Uncovering new families and folds in the natural protein universe. *Nature.*  
757 2023;622(7983).
- 758 34. Del Alamo D, Sala D, McHaourab HS, Meiler J. TITLE: Sampling alternative  
759 conformational states of transporters and receptors with AlphaFold2. *Elife.* 2022;11.
- 760 35. Wayment-Steele HK, Ojoawo A, Otten R, Apitz JM, Pitsawong W, Hömberger M, et  
761 al. Predicting multiple conformations via sequence clustering and AlphaFold2. *Nature.*  
762 2024;625(7996).
- 763
- 764
- 765
- 766
- 767
- 768
- 769
- 770
- 771
- 772
- 773
- 774
- 775
- 776

## 777 7. Tables & Figures

### 778 7.1 Tables

<b>Method</b>	<b>Webserver</b>	<b>Queries</b>	<b>Databases</b>
A	clusters.foldseek.com	AbnU, AbyU, Cyc15,Pyrl4	AlphaFold/Uniprot50 v4 AlphaFold/Swiss-prot v4 AlphaFold/Proteome v4 PDB100
B	search.foldseek.com	AbnU, AbyU, Cyc15,Pyrl4	AlphaFold/Proteome v4
C	<a href="#">search.COFACTOR</a>	AbyU	PDB100
D	<a href="#">search.plmBLAST</a>	AbnU, AbyU, Cyc15,Pyrl4	PDB100

779 **TABLE 1** - Detailed Search Protocol: Queries, Methods, and Databases

780 Detailed search methodologies, including the specific search queries, methods employed,  
 781 and databases utilized for data collection. The table provides a comprehensive overview of  
 782 the approach used in each search instance, outlining the search parameters, the databases  
 783 accessed, and the strategies applied to gather relevant data for the study. This structured  
 784 format ensures transparency and reproducibility of the search process.

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

	<b>AbnU</b>	<b>AbyU</b>	<b>Cyc15</b>	<b>PyrI4</b>
<b>AbnU</b>	100%	49%	31%	<sup>801</sup> 26% <sub>803</sub>
<b>AbyU</b>	49%	100%	39%	<sup>802</sup> 31% <sub>804</sub>
<b>Cyc15</b>	31%	39%	100%	<sup>805</sup> 31% <sub>806</sub>
<b>PyrI4</b>	26%	31%	26%	100% <sub>807</sub>
<b>A0A246RFQ9</b>	83%			<sup>808</sup>

809 **TABLE 2-** Analysis of Sequence Similarity Among Queries: Insights and Comparisons.  
 810 Sequence similarity analysis of queries, including A0A246RFQ9, to assess the relationship  
 811 between sequence similarity and the retrieval process. A0A246RFQ9 was selected based on  
 812 a multiple sequence alignment (MSA) performed with ClustalW, where it showed the highest  
 813 sequence similarity to AbyU. Sequence similarity was calculated using BLASTp to evaluate if  
 814 higher similarity between queries influences their results. The table displays the similarity  
 815 scores between the sequences, providing insight into the comparison of query sequences  
 816 used in the study.

817  
 818  
 819  
 820  
 821  
 822  
 823  
 824  
 825  
 826  
 827  
 828  
 829  
 830  
 831  
 832  
 833  
 834  
 835  
 836

Subject			
Query	AbyU	AbnU	A0A246RFQ9
AbyU	0	1123	408
AbnU	1968	0	
A0A246RFQ9	1274	0	0
Total hits	2345	1500	1479

**TABLE 3 - Relationship Between Sequence Similarity and Unique Hit Retrieval.**

837 Unique hit analysis for query A0A246RFQ9 using Foldseek. A0A246RFQ9 was included due  
 838 to its high sequence similarity with AbyU. The table compares the number of unique hits  
 839 retrieved by A0A246RFQ9 and the subject, allowing for an assessment of how the query's  
 840 high sequence similarity influences the retrieval of unique hits. This analysis provides insight  
 841 into the relationship between sequence similarity and the distinctiveness of the hits retrieved.  
 842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

<b>Query</b>	<b>Per query</b>	<b>Only 3Di</b>	<b>Only TM-align</b>	<b>Total</b>
AbnU	694	720	780	1500
AbyU	1357	1126	1219	2345
Cyc15	2054	1463	1585	3048
Pyrl4	1538	1130	1224	2354

867 **TABLE 4 – Unique Hit Distribution Across Queries: Foldseek Method A Analysis.**

868 Unique hits retrieved by each alignment or query using Foldseek (Method A). Unique hits are  
 869 defined as those retrieved exclusively by individual alignments or queries. These hits were  
 870 obtained by merging all the results from Foldseek and comparing the hits retrieved by each  
 871 alignment and query. The table highlights the distinct hits retrieved by each query,  
 872 showcasing the unique retrieval profiles based on the Foldseek analysis.

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

<b>Query</b>	<b>Cluster representative</b>	<b>Total members</b>
AbnU	A0A0K9XAS6	78
AbyU	A0A0K9XAS6	78
Additional	A0A0B6VM83	37
Additional	A0A022Q700	4092
Additional	A0A2P1BT29	2
Additional	A0A4Y8Y1N5	3
Additional	A0A4P7ZYP8	2
Additional	A0A833W1H9	49
Additional	L8H095	126
Cyc15	A0A0K9XAS6	78
Pyrl4	A0A3A9ZPZ1	38
<b>Total hits</b>		<b>4583</b>

899 **TABLE 5 – Query Distribution Across Clusters: Method B and Search-Derived Additions.**  
900 Cluster analysis of queries using Method B. In this method, each query is assigned to a  
901 cluster, with each cluster represented by a designated representative sequence. The table  
902 presents the number of members in each cluster, reflecting the distribution of queries within  
903 their respective clusters. Additional clusters identified during the query searches were  
904 included in the analysis, providing a comprehensive overview of the cluster memberships  
905 and the additional clusters discovered during the process.  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923

<b>Method</b>	<b>Total hits (relevant)</b>	<b>Unique hits</b>	<b>Positive control</b>
Method A	120	76	AbmU, AbsU, AbnU, Cyc15, AbyU, ChIL, TcaU4, LonU2, VstJ, LobD1, Pyrl4, Tmn8, Tsn15
Method B	45	5	AbmU, AbnU, AbyU, Cyc15, QmnH, ChIL, TcaU4, Tmn8, Tsn15

**TABLE 6- Unique Hit Retrieval and Positive Control Detection: Method A vs. Method B.**

924 Comparison of unique hits obtained by both Method A (Foldseek) and Method B. The table  
 925 shows the count of unique hits retrieved by each method. Additionally, it includes an  
 926 assessment of the presence of positive controls within the results, providing insight into the  
 927 effectiveness and accuracy of each method in identifying unique hits and detecting positive  
 928 controls.

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

Search type	Unique hits	Total hits
Sequence-based search	36	76
Structure-based search	58	98

**TABLE 7-** Comparative Effectiveness: Structure vs. Sequence in Hit Retrieval.

Analysis of hits obtained from Zorn *et al.*, (sequence hits) with variant hits excluded. The table presents the count of unique hits retrieved by each method, after discarding any variant hits reported by Zorn et al. A comparative analysis of the unique hits retrieved by each method is provided, offering insights into the consistency and effectiveness of the methods in identifying distinct hits. This table effectively compared the two different approaches of using sequence or structure for doing a homolog search.

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

<b>Uniprot ID</b>	<b>Salt bridge (78 + 122 AbyU)</b>
A0A0B6VRF8	A-R
A0A0F7N0P6	Y-L
A0A0G3A932	E-R
A0A0J1H8G7	E-R
A0A0K9XAS6	E-R
A0A0L8NIA4	E-R
A0A0M2RZN8	N-R
A0A0M8XYC3	E-R
A0A0Y0AWW3	E-R
A0A101JBC1	E-R
A0A109B5R0	E-R
A0A132CHM2	E-R
A0A168IKI6	E-R
A0A1B4FKD2	E-R
A0A1B4FUY1	E-R
A0A1B4SIL4	E-R
A0A1C3HDP6	E-R
A0A1C4PF23	E-E
A0A1C6V0W8	T-M
A0A1J4Q2U2	A-R
A0A1M7FPM3	Y-R
A0A1Q4P580	Y-R
A0A1Q5TWR6	Y-R
A0A1Q7WCG8	D-Y
A0A1S1QVM7	E-R
A0A1S1R073	Y-R
A0A1S1R326	E-R
A0A1S2PAK5	L-M
A0A1Y0RPR1	L-I
A0A1Y2SP38	Y-R
A0A239CWW2	E-R
A0A246RFQ9	E-R
A0A248YPZ6	I-M
A0A2P1BT29	Y-R
A0A2P4UBX2	I-W
A0A2S6GLP7	Y-R
A0A2S6GLP9	R-R
A0A2S6H1S2	F-R
A0A318JTK0	D-T
A0A345XXX2	Y-R
A0A367F712	S-M
A0A3A9YQ06	I-R
A0A3A9ZPZ1	S-E
A0A3D1PI25	E-R

---

A0A3E0H6P9	F-R
A0A401MXE6	V-R
A0A401MXH6	V-R
A0A438LZ58	T-E
A0A495Q9P3	Y-R
A0A4R2JL81	Y-R
A0A4R4MHI3	
A0A4R5A4Q9	Y-R
A0A4R5E5X7	T-L
A0A4U3F911	Y-R
A0A4U3LVR5	Y-R
A0A4Y8Y1N5	Y-R
A0A518WFP6	E-R
A0A544YAL3	E-R
A0A5D0NLL8	A-R
A0A5P8WEY5	E-I
A0A5P9PSM5	L-M
A0A5Q0GZD2	L-M
A0A5R8YLG6	F-R
A0A5S4GQR6	F-R
A0A5S4H4D6	V-R
A0A5S8WF63	V-L
A0A640SQ86	E-R
A0A640SVV6	E-R
A0A6G3UFY1	Y-R
A0A6G9F6I2	L-M
A0A6I4MT57	S-L
A0A6N7ZAT6	F-R
A0A6N7ZB01	S-R
A0A743Z1Y4	E-R
A0A7D3ZDF8	V-L
A0A7J5DHM5	E-R
A0A7K2LH80	E-R
A0A7K2N5N2	E-M
A0A7U4PB97	E-R
A0A7W4FA88	E-R
A0A7X0IJX5	E-R
A0A7Y0J7D7	E-R
A0A7Y0WQ23	E-R
A0A8B4GMM4	E-R
A7BEY9	S-L
F4F7G1	E-R
F8AWP7	E-R
F8K320	N-M
K4FDH3_2	
K4FDH3_1	E-R

---

---

K7QVW7	S-
L7RSB5	A-R
Q0R4M0	S-M

987 **TABLE 8-** Comparative Analysis of Salt Bridges at AbyU-Equivalent 78-122 Positions.  
988 Evaluation of the salt bridge composition between Glu (78) and Arg (122) AbyU positions.  
989 The presence of this salt bridge was assessed through structural alignment of the hits. The  
990 table details the alignment results and the subsequent visualization of the structure in  
991 PyMOL to confirm the position of the salt bridge. This analysis provides insight into the  
992 structural features of the hits and their relevance to the salt bridge configuration compared to  
993 that of AbyU.

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

<b>UniProt ID</b>	<b>Conformation in AFDB</b>	<b>Conformation in PDB</b>
A0A0B6VRF8	Open	Not available
A0A0F7N0P6	Open	Not available
A0A0G3A932	Closed	Not available
A0A0J1H8G7	Closed	Not available
A0A0K9XAS6	Closed	Not available
A0A0L8NIA4	Closed	Not available
A0A0M2RZN8	Open	Not available
A0A0M8XYC3	Closed	Not available
A0A0Y0AWW3	Closed	Not available
A0A101JBC1	Closed	Not available
A0A109B5R0	Closed	Not available
A0A132CHM2	Closed	Not available
A0A168IKI6	Closed	Not available
A0A1B4FKD2	Closed	Not available
A0A1B4FUY1	Closed	Not available
A0A1B4SIL4	Closed	Not available
A0A1C3HDP6	Closed	Not available
A0A1C4PF23	Open	Not available
A0A1C6V0W8	Open	Not available
A0A1J4Q2U2	Open	Not available
A0A1M7FPM3	Open	Not available
A0A1Q4P580	Closed	Not available
A0A1Q5TWR6	Closed	Not available
A0A1Q7WCG8	Open	Not available
A0A1S1QVM7	Closed	Not available
A0A1S1R073	Closed	Closed
A0A1S1R326	Closed	Not available
A0A1S2PAK5	Closed	Not available
A0A1Y0RPR1	Open	Not available
A0A1Y2SP38	Closed	Not available
A0A239CWW2	Closed	Not available
A0A246RFQ9	Closed	Not available
A0A248YPZ6	Open	Not available
A0A2P1BT29	Open	Open
A0A2P4UBX2	Open	Not available

A0A2S6GLP7	Closed	Not available
A0A2S6GLP9	Closed	Not available
A0A2S6H1S2	Closed	Not available
A0A318JTK0	Loop not defined	Not available
A0A345XXX2	Open	Not available
A0A367F712	Open	Not available
A0A3A9YQ06	Loop not defined	Not available
A0A3A9ZPZ1	Open	Not available
A0A3D1PI25	Loop not defined	Not available
A0A3E0H6P9	Closed	Not available
A0A401MXE6	Closed	Open
A0A401MXH6	Closed	Not available
A0A438LZ58	Open	Not available
A0A495Q9P3	Closed	Not available
A0A4R2JL81	Closed	Not available
A0A4R4MHI3	Closed	Not available
A0A4R5A4Q9	Open	Not available
A0A4R5E5X7	Open	Not available
A0A4U3F911	Closed	Not available
A0A4U3LVR5	Closed	Not available
A0A4Y8Y1N5	Open	Not available
A0A518WFP6	Closed	Not available
A0A522WLH4	Non-beta barrel fold	Not available
A0A544YAL3	Closed	Not available
A0A5C8PZY7	Open	Not available
A0A5D0NLL8	Closed	Not available
A0A5P8WEY5	Loop not defined	Not available
A0A5P9PSM5	Closed	Not available
A0A5Q0GZD2	Closed	Not available
A0A5R8YLG6	Closed	Not available
A0A5S4GQR6	Closed	Not available
A0A5S4H4D6	Open	Not available
A0A5S8WF63	Open	Not available
A0A640SQ86	Open	Not available
A0A640SVV6	Closed	Not available
A0A6G3UFY1	Open	Not available

A0A6G9F6I2	Open	Not available
A0A6I4MT57	Open	Not available
A0A6N7ZAT6	Open	Not available
A0A6N7ZB01	Closed	Not available
A0A743Z1Y4	Closed	Not available
A0A7D3ZDF8	Open	Not available
A0A7J5DHM5	Open	Not available
A0A7K2LH80	Closed	Not available
A0A7K2N5N2	Closed	Not available
A0A7U4PB97	Closed	Not available
A0A7W4FA88	Closed	Not available
A0A7X0IJX5	Closed	Not available
A0A7Y0J7D7	Closed	Not available
A0A7Y0WQ23	Closed	Not available
A0A838IMS4	Non-beta barrel fold	Not available
A0A8B4GMM4	Closed	Not available
A7BEY9	Open	Not available
F4F7G1	Closed	closed and open
F8AWP7	Closed	Not available
F8K320	Closed	Not available
K4FDH3_1	Open	Not available
K4FDH3_2	Open	Not available
K7QVW7	Open	closed and open
L7RSB5	Open	Not available
Q0R4M0	Open	Not available

1020 **TABLE 9-** Evaluation of Hit Conformations and Their Representation in the PDB.

1021 Assessment of structural conformations and availability in the PDB. The table categorizes  
 1022 hits based on their structural conformations: some exhibited a closed conformation, while  
 1023 others lacked a defined capping loop, making their conformation indeterminate. Additionally,  
 1024 the presence of these structures in the Protein Data Bank (PDB) was evaluated, providing  
 1025 insights into the structural characteristics and reported conformations.

1026

1027

1028

1029

1030

UniProt ID	Template	Result
A0A0F7N0P6	AbyU	Open
A0A0F7N0P6	Cyc15	Open
A0A0M2RZN8	AbyU	Open
A0A0M2RZN8	AbyU	Open
A0A1C6V0W8	AbyU	Open
A0A1M7FPM3	AbyU	Open
A0A1M7FPM3	AbyU	Open
A0A1Q7WCG8	AbyU	Open
A0A1S2PAK5	AbyU	Open
A0A1S2PAK5	AbyU	Open
A0A1Y0RPR1	AbyU	Open
A0A1Y0RPR1	Cyc15	Open
A0A248YPZ6	AbyU	Open
A0A2P4UBX2	AbyU	Open
A0A345XXX2	AbyU	Closed
A0A367F712	AbyU	Open
A0A3A9ZPZ1	AbyU	Open
A0A438LZ58	AbyU	Open
A0A4R5E5X7	AbyU	Open
A0A4Y8Y1N5	AbyU	Open
A0A5P8WEY5	Cyc15	Open
A0A5Q0GZD2	AbyU	Open
A0A5S8WF63	AbyU	Open
A0A6G9F6I2	AbyU	Open
A0A6I4MT57	AbyU	Open
A0A6I4MT57	AbyU	Open
A0A7D3ZDF8	AbyU	Open
A0A7J5DHM5	AbyU	Open
A0A7K2N5N2	AbyU	Open
A7BEY9	AbyU	Open
F8K320	AbyU	Open

1031 **TABLE 10-** Template-Based Closed Conformation Analysis: AbyU and Cyc15 in ColabFold.  
 1032 Closed conformation prediction of structural hits performed with ColabFold. AbyU and Cyc15  
 1033 were used as templates for the predictions, selected based on their sequence similarity to  
 1034 the subject (refer to Figure SX for details). This table summarizes the predictions made  
 1035 using these templates to assess the closed conformations of the structural hits.

Test	Model	Subject (UniProt ID)	Template	Num_seeds	drop_out	max_ms	Result
Test_1_1	Model 1	A0A1Y0RPR1	AbyU	2	checked	256:51	no change
Test_1_2	Model 2	A0A1Y0RPR1	AbyU	2	checked	512:10	no change
Test_1_3	Model 1	A0A1Y0RPR1	AbyU	2	checked	64:128	no change
Test_1_4	Model 2	A0A1Y0RPR1	AbyU	2	checked	32:64	no change
Test_1_5	Model 2	A0A1Y0RPR1	AbyU	2	checked	16:32	no change
Test_2_1	Model 1	A0A1Y0RPR1	AbyU	16	checked	256:51	alternative conformation
Test_2_2	Model 1	A0A1Y0RPR1	AbyU	16	checked	16:32	alternative conformation
Test_2_3	Model 1	A0A1Y0RPR1	AbyU	16	checked	512:10	alternative conformation
Test_2_4	Model 2	A0A1Y0RPR1	AbyU	16	checked	64:128	alternative conformation
Test_2_5	Model 3	A0A1Y0RPR1	AbyU	16	checked	32:64	alternative conformation
Test_3_1	Model 1	A0A1Y0RPR1	AbyU	2	checked	default	no change
Test_3_1	Model 2	A0A1Y0RPR1	AbyU	2	checked	default	alternative conformation
Test_3_1	Model 3	A0A1Y0RPR1	AbyU	2	checked	default	alternative conformation
Test_3_1	Model 4	A0A1Y0RPR1	AbyU	2	checked	default	no change
Test_4_1	Model 1	A0A1Y0RPR1	AbyU	4	checked	default	no change
Test_4_1	Model 2	A0A1Y0RPR1	AbyU	4	checked	default	alternative conformation
Test_4_1	Model 3	A0A1Y0RPR1	AbyU	4	checked	default	alternative conformation
Test_4_1	Model 4	A0A1Y0RPR1	AbyU	4	checked	default	alternative conformation
Test_5_1	Model 1	A0A1Y0RPR1	AbyU	8	checked	default	no change
Test_5_1	Model 2	A0A1Y0RPR1	AbyU	8	checked	default	no change
Test_5_1	Model 3	A0A1Y0RPR1	AbyU	8	checked	default	alternative conformation
Test_5_1	Model 4	A0A1Y0RPR1	AbyU	8	checked	default	alternative conformation
Test_6_1	Model 1	A0A1Y0RPR1	AbyU	16	checked	default	no change
Test_6_1	Model 2	A0A1Y0RPR1	AbyU	16	checked	default	no change
Test_6_1	Model 3	A0A1Y0RPR1	AbyU	16	checked	default	no change
Test_6_1	Model 4	A0A1Y0RPR1	AbyU	16	checked	default	no change

1036 **TABLE 11-** Optimized ColabFold Predictions: Assessing Structural Hit Conformations.  
1037 Closed conformation prediction of structural hits using ColabFold. The predictions were  
1038 made with adjustments to the multiple sequence alignment (MSA) depth and number of  
1039 seeds, based on prior recommendations. For Test 1 and Test 2, the highest-scoring models  
1040 were selected for visualization. In Tests 4, 5, and 6, the four highest-scoring models were  
1041 visualized. A conformation was classified as alternative if the capping loop was positioned  
1042 differently compared to the original subject.

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

	<b>Total</b>	<b>Mean volume</b>
Hits in a closed conformation	54	
Hits in a closed conformation with the salt bridge	51	124
Hits in a closed conformation without the salt bridge	3	229
Hits in an open conformation	34	
Total hits	96	

1073 **TABLE 12-** Comparative Cavity Volume Assessment of Closed Conformation Hits Using

1074 CastP.

1075 Calculation of cavity volumes for hits in a closed conformation using CastP. The table  
 1076 includes the volume measurements for hits with a closed conformation, with separate  
 1077 calculations provided for hits lacking the salt bridge (or without Arg at position 122). This  
 1078 distinction allows for the evaluation of how the presence or absence of the salt bridge affects  
 1079 the cavity volumes of the hits.

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

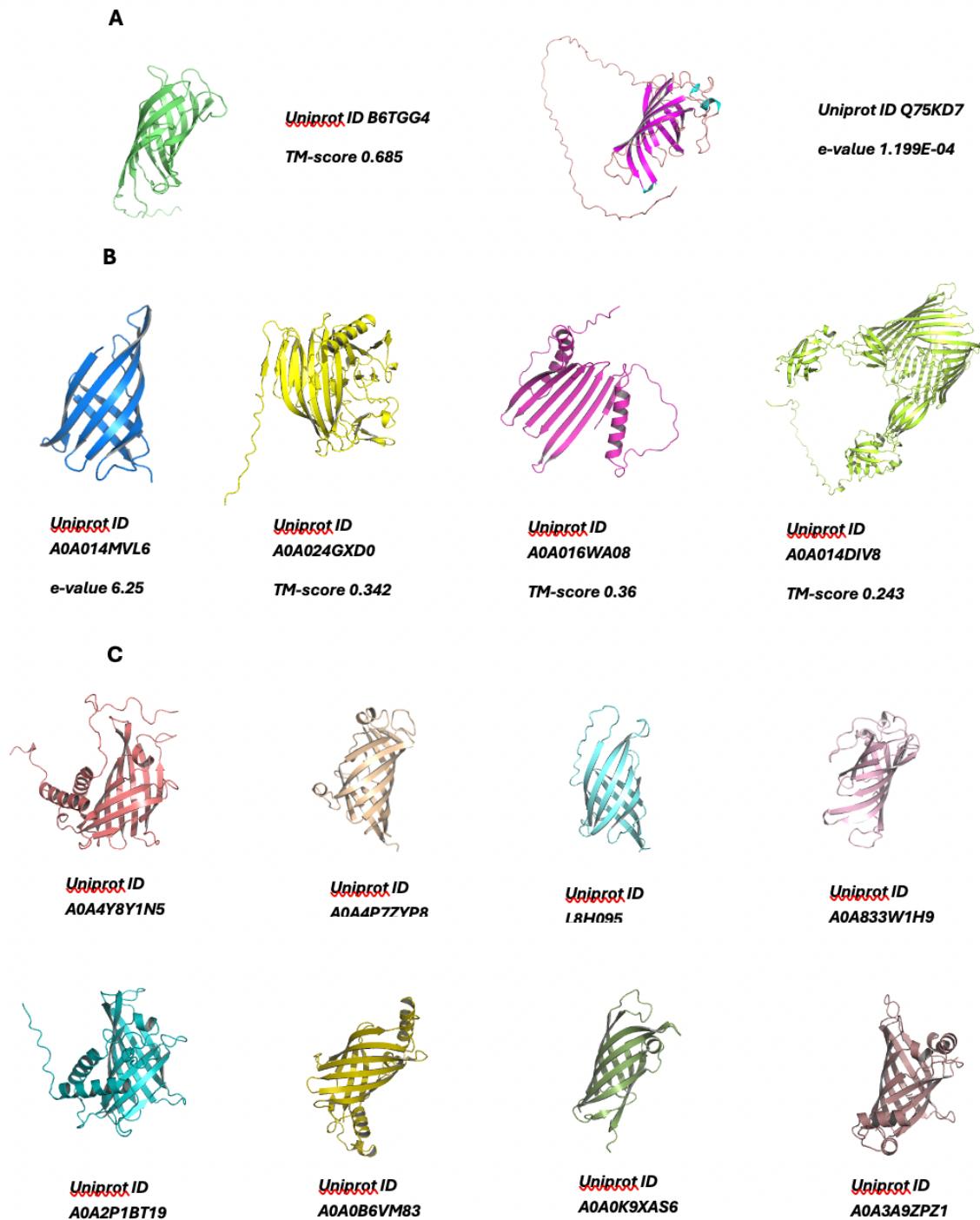
1101

1102

1103

1104 6.2 Figures

1105



1106

1107 **FIGURE 1-** Overview of unique hits and cluster representatives retrieved by alignment  
1108 methods.

1109 (A): Unique hits retrieved by the alignment method. These hits display similar structures,  
1110 indicating no particular retrieval advantage for each alignment method. Additionally, these  
1111 hits are considered irrelevant to the study.

1112 (B): Unique hits retrieved by individual queries. These hits also exhibit non-beta barrel folds  
1113 or incomplete beta barrels, rendering them irrelevant. This finding supports the accuracy of  
1114 information density capture by Foldseek, as these hits have lower e-values and TM scores,  
1115 demonstrating that Foldseek correctly identifies relevant structures.

1116 (C): Cluster representatives of each cluster retrieved by Method B. The representatives  
1117 display similar structures to the queries, even though some are uncharacterized proteins.  
1118 This indicates that Method B effectively clusters proteins with similar structural features.

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

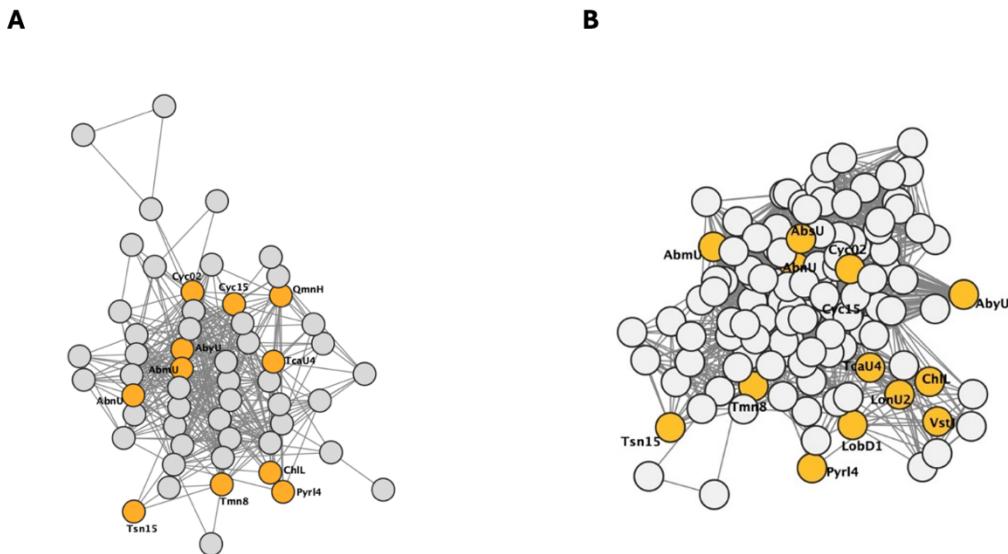
1135

1136

1137

1138

1139



1140

1141 **FIGURE 2-** Comparison of Sequence Similarity Networks (SSNs) for Method A and Method  
1142 B.

1143 **(A):** SSN for Method A (Foldseek), displaying all positive controls clustered with other hits in  
1144 a separate cluster. The positive controls identified in this network include: AbmU, AbsU,  
1145 AbnU, Cyc15, AbyU, ChIL, TcaU4, LonU2, VstJ, LobD1, Pyrl4, Tmn8, and Tsn15. Hits within  
1146 this cluster are considered relevant as they are grouped with the positive controls.

1147 **(B):** SSN for Method B (AFDB Clusters), showing all positive controls clustered with other  
1148 hits in a separate cluster. The positive controls in this network are: AbmU, AbnU, AbyU,  
1149 Cyc15, QmnH, ChIL, TcaU4, Tmn8, and Tsn15. Hits within this cluster are considered  
1150 relevant as they are grouped with the positive controls.

1151

1152

1153

1154

1155

1156

1157

1158

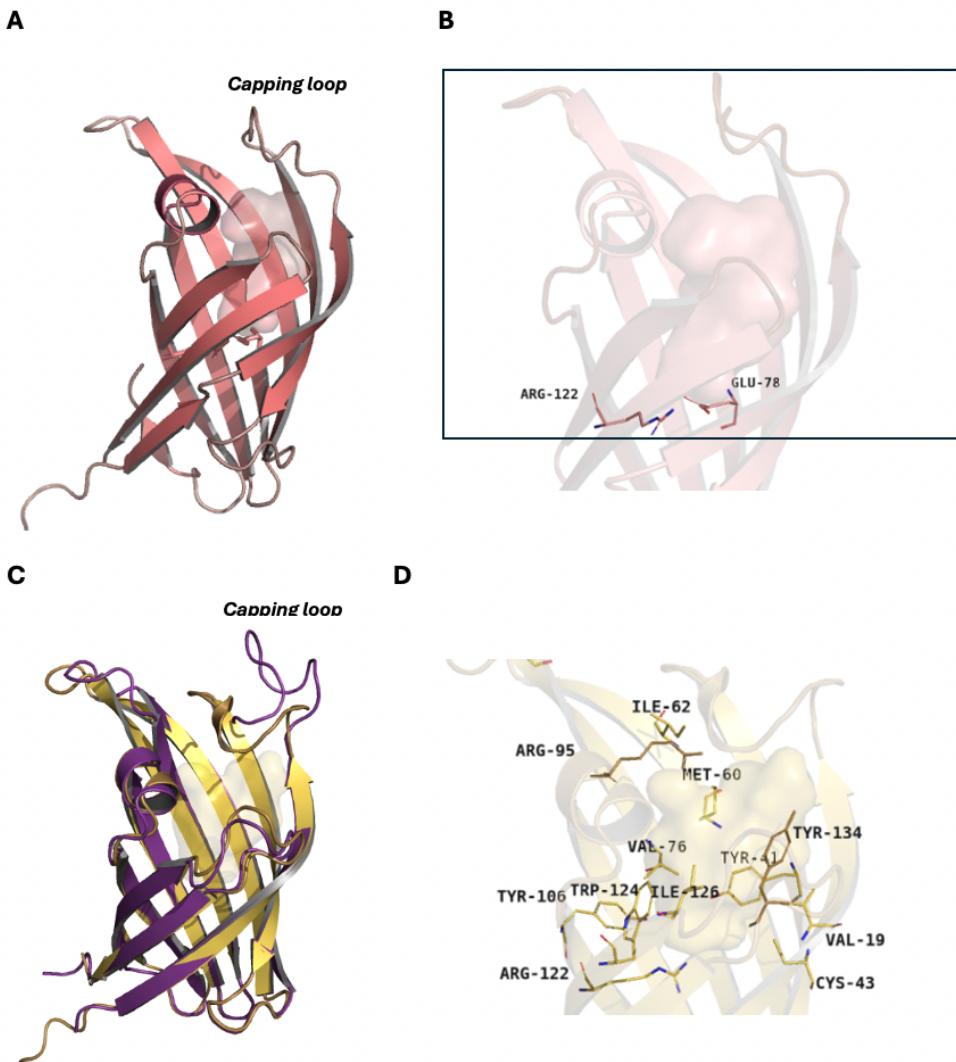
1159

1160

1161

1162

1163



1164

1165 **FIGURE 3-** Structural visualization of AbyU and Cyc15 with emphasis on cavities and salt  
1166 bridges.

1167 (A): Complete structure of AbyU, highlighting the cavity delimited by the salt bridge between  
1168 Glu 78 and Arg 122. The salt bridge is represented by sticks, and the position of the capping  
1169 loop is indicated.

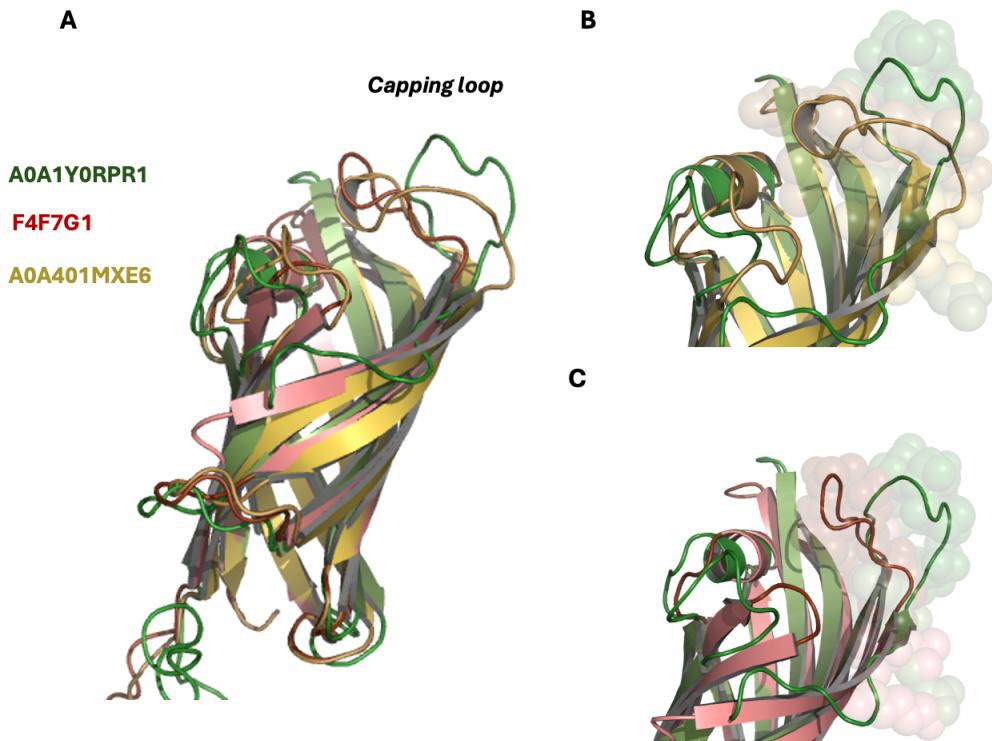
1170 (B): Close-up view of the AbyU structure, focusing on the surface representation of the cavity  
1171 and the salt bridge. This detailed view provides a closer look at the cavity's dimensions and  
1172 the salt bridge.

1173 (C): Comparison of Cyc15 in its open conformation (PDB 8OF7) and closed conformation  
1174 (Alpha Fold prediction based on template UniProt A0A0K9XAS6). The closed conformation  
1175 is shown in yellow and the open conformation in purple.

1176 (D): Detailed close-up of Cyc15, highlighting the cavity-lining residues that define the cavity  
1177 in the closed conformation.

1178

1179



1180

1181 **FIGURE 4-** Closed conformation predictions using Cyc15 and AbyU.1182 (A): Alignment of A0A1Y0RPR1 with AbyU and Cyc15, highlighting the capping loop in  
1183 sphere representation. This panel provides a comparative overview of the alignment and the  
1184 positioning of the capping loop when using each template.1185 (B): Close-up view of the capping loop of A0A1Y0RPR1, showing that its conformation did  
1186 not change when Cyc15 was used as a template. This panel demonstrates the stability of  
1187 the capping loop conformation with the Cyc15 template.1188 (C): Close-up view of the capping loop of A0A1Y0RPR1, showing that its conformation did  
1189 not change when AbyU was used as a template. This panel highlights the consistency of the  
1190 capping loop conformation with the AbyU template.1191 The predictions with varying parameters were conducted after these attempts were  
1192 unsuccessful, indicating the need for further adjustments.

1193

1194

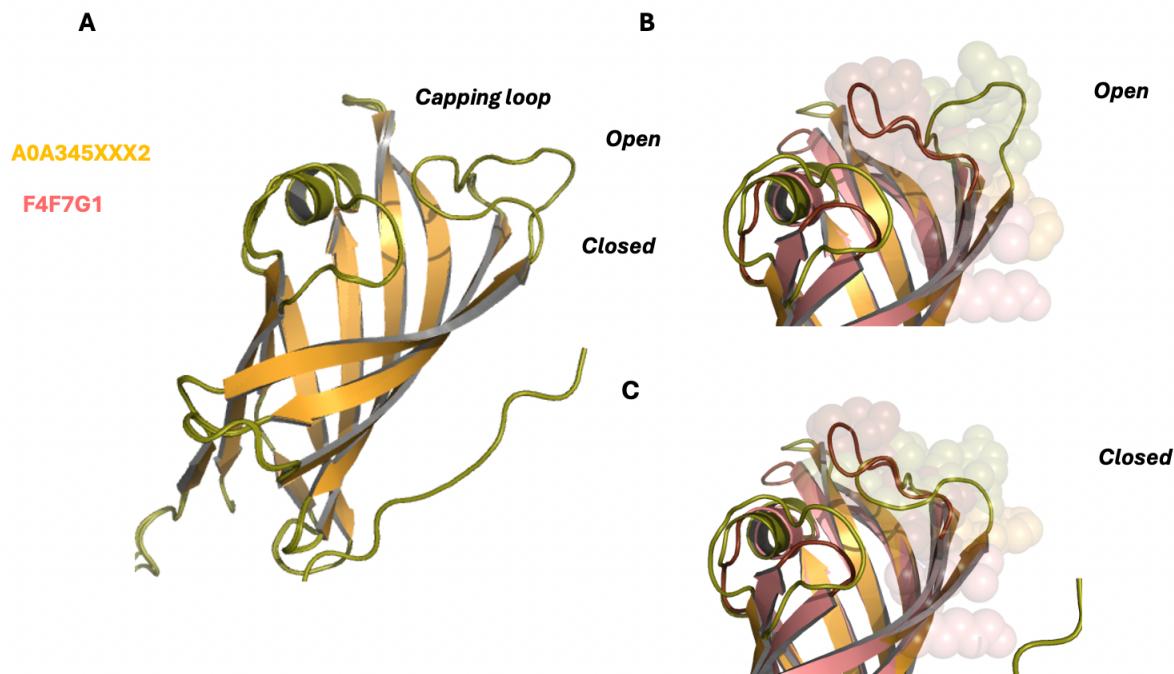
1195

1196

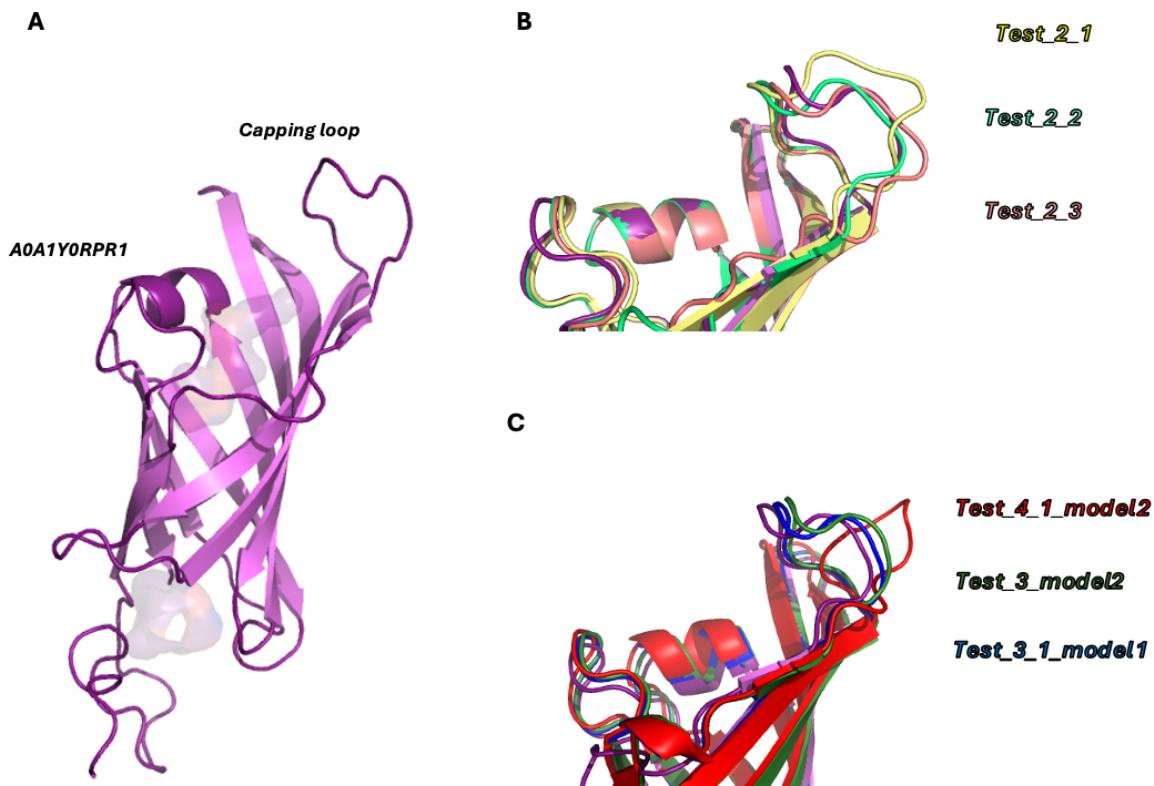
1197

1198

1199



1200  
1201 **FIGURE 5-** Successful closed conformation prediction of A0A345XXX2 using AbyU as a  
1202 template.  
1203 (A): Structure of A0A345XXX2 aligned with AbyU, highlighting the capping loop. Initially,  
1204 A0A345XXX2 had an open loop conformation. Using AbyU as a template with default  
1205 parameters in ColabFold resulted in a successfully predicted closed conformation.  
1206 (B): Close-up view showing the open conformation of the capping loop before the alignment  
1207 with AbyU. This panel illustrates the initial structure of the capping loop.  
1208 (C): Close-up view of the closed conformation obtained for A0A345XXX2 using AbyU as a  
1209 template. This panel highlights the successfully predicted closed state of the capping loop.  
1210  
1211  
1212  
1213  
1214



1215

1216

1217 **FIGURE 6-** Prediction of the closed conformation of A0A1Y0RPR1 through variations in  
1218 MSA depth and number of seeds.

1219 (A): Structure of A0A1Y0RPR1 depicted in pink, with its cavity shown in surface  
1220 representation. This view provides a general overview of the protein structure and its internal  
1221 cavity.

1222 (B): Close-up of the capping loop of A0A1Y0RPR1, aligned with the results from Test\_2\_1,  
1223 Test\_2\_2, and Test\_2\_3. This panel illustrates the different conformations of the capping  
1224 loop resulting from variations in MSA depth and the number of seeds used in the predictions.

1225 (C): Results from Test\_4\_1\_Model2, Test\_3\_Model2, and Test\_3\_1\_Model1, where only the  
1226 number of seeds was varied while the MSA depth remained constant. This panel  
1227 demonstrates that the capping loop did not acquire significantly different conformations  
1228 compared to the results in panel (B), highlighting the relative stability of the loop  
1229 conformation with changes in the number of seeds.

1230

1231

1232

1233

1234

1235 **Appendix**

1236 *Supplementary Tables*

<b>Query Uniprot ID</b>	<b>Subject Uniprot ID</b>	<b>Identity (%)</b>	<b>Coverage (%)</b>
A0A0F7N0P6	A0A3S8VII0	100	100
	A0A2K9FB81	100	100
A0A0M2RZN8	A0A0B6VM83	100	100
A0A6G3UFY1	A0A1H5HAZ2	100	100
	A0A6G2Z021	100	100
	A0A6I7XZ82	100	100
	A0A327URC2	100	100
	A0A6B3BH80	99	100
	A0A4S2UVU8	97	100
B5L6K7	A0A1S8XP27	99	100
A0A7Y0J7D7	A0A6G3VZ74	100	100
	A0A6I5GVW7	100	100
	A0A7K2LZ25	100	100
	A0A1V0QH49	100	100
	D6K898	99	100
F4F7G1	A0A097CSY6	100	100
M9T3W4	A0A5P9NYB6	100	100
	L7RSB5	97	100
	A0A6G9ENB7	96	100
A0A0L8NIA4	A0A0N0T4L8	100	100
A0A5S8WF63	A0A5C8PZY7	100	99
	A0A6C0QDS2	100	99
A0A5R8YLG6	A0A5D0QA99	100	100
A0A1B4FKD2	A0A7T2X384	100	96
L8EX63	A0A7K2RBB9	100	100
	A0A429JRF8	100	100
	A0A429I7L3	100	100
A0A5D0NLL8	A0A5D0UIN8	99	100
	A0A5D3FXF0	99	100
A0A2P1BT29	A0A4Q6VF37	99	100
A0A4U3LVR5	A0A7C9NN96	100	100
A0A1S1QVM7	A0A166S4A2	97	100
A0A2P4UBX2	A0A6I4WC81	97	100
A0A367F712	A0A367ECF2	97	100
A0A4R4MHI3	A0A7W7IEZ4	97	100
	A0A4R4TSZ3	97	100
A0A1C4PF23	A0A7K2G773	100	100
A0A0M8XYC3	A0A7K3G8J5	96	100
A0A4R5E5X7	A0A4R4TSZ3	96	100

1237 **TABLE S1-** Variants in structural hits found by BlastP

1238 This table presents an analysis of structural variants found in protein sequences identified as  
1239 hits through BLASTP searches.

1240

1241

1242

1243

Cyclase name	Unitprot ID	Reference
AbyU	F4F7G1	
AbmU	A0A2P1BT29	(Li et al., 2022)
AbnU	A0A1S1R073	(Kashyap et al., 2021)
AbsU	<a href="#">A0A1V0QH49</a>	(Xu and Yang, 2021)
Cyc15	A0A401MXE6	(Zorn et al., 2023)
Pyrl4	K7QVW7	(Tian et al., 2015)
LobD1	L7RSB5	(Yue et al., 2016)
Tsn15	A0A5S8WF63	(Fujiyama et al., 2021)
TcaU4	B5L6K7	(Fang et al., 2008)
VstJ	A0A0B6VRF8	(Hashimoto et al., 2015)
Tmn8	<a href="#">A7BEY9</a>	(Demydchuk et al., 2008)
ChIL	Q0R4M0	(Jia et al., 2006)
QmnH	K4FDH3	(Jeon et al., 2017)

1244 **TABLE S2-** Positive control

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

Cyclase name	Uniprot ID	Fasta
AbnU	A0A1S1R073	MSVVHEGIWEPQIRNEQNVNVADPQVGQIGSYYDEL YDSSRELLGITIGRYEIRYKKVGGAVLTYYSEDLFLRD GIIHAEGWADFNDVKNGVWVGYPAVGLDGVRGLD GRREWRVIEPDQPVEARISLHG
AbyU	F4F7G1	MTERLETRPQALLIKVPTEIVVKVWDDVDVAAPAVGQ VGKFDDELYDEAGAQIGTSSGNFRIEYVRPTDGGLLT YYQEDITLSDGVIAEGWADFNDVRTSKWVFYPATG VSGRYLGLTGFRQWRMTGVRKSAEARILLGE
Cyc15	A0A401MXE6	MGDTTTANDVVTVELVEKVTKKDLNESGSIEGFGPG MMATYWCDVFDTEGKHIGTTVGCMIDILYADPESGHL VEHVAEQIRLPDGTIMAWGTMNRSDVLAQKWITYRC QGTSGRYAGLVGTRTWRIQSLEDESYPIVAKMELRG A
Pyrl4	K7QVW7	MTTPQIDERAMEAGAAALQETIVDPGPLDVTALAVAA ALAAGLHSAADDPAAALDKCIVLDELTEFAEKLVVHD RPGGIGTTVEYVEVYEDASGVRLGTATGNAVVLKME PHMWQFHQSSELADGSFEAVGVIDCTAMLRRMTQ VLRVTGRSGRYAGKSGFMTLAISDPNQRPPHYSVQV VLC

**TABLE S3-** fasta sequences of queries AbnU, AbyU, Cyc15 and Pyrl4

1270  
 1271  
 1272  
 1273  
 1274  
 1275  
 1276  
 1277  
 1278  
 1279  
 1280  
 1281  
 1282  
 1283  
 1284  
 1285  
 1286  
 1287  
 1288  
 1289  
 1290

Method	Minimum length of sequence	Maximum length of sequence	Alignment score threshold
Method A	120	625	7
Method B	120	625	7

**TABLE S5-** Parameters used for Sequence similarity Network visualization.

1291

1292

1293

1294

1295

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

Query	Total hits per query
AbnU	23
AbyU	21
Cyc15	23
Pyrl4	18
<b>Total hits</b>	<b>85</b>

1325 **TABLE S6-** Hits retrieved by method C

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

1350

1351

1352

1353

1354

<b>Query</b>	<b>Total hits</b>
AbyU	10
<b>List of hits</b>	
<b>Uniprot ID</b>	<b>PDB</b>
A0A1S1R073	7DVI
A0A1V1FH01	6OOC
A0A2P1BT29	6LE0
A0A401MXE6	8OF7
A0A5S8WF63	6NNW
A0A9Y2YAD1	7X7Z
F4F7G1	5DYQ
K0V2D8	5K21
K7QVW7	5BTU
Q8H0N6	4H69

1355 **TABLE S7-** Hits retrieved by method D

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

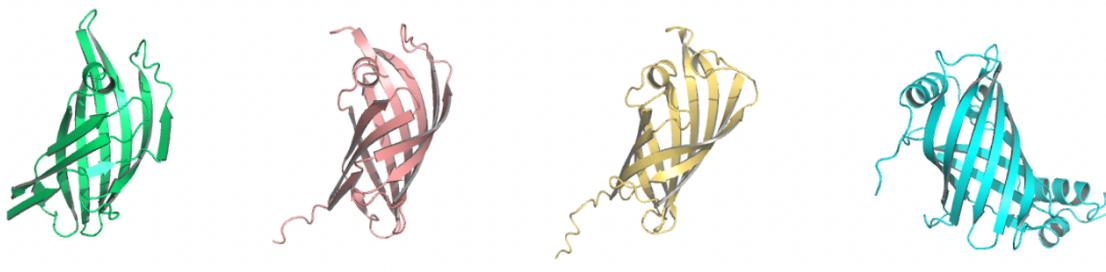
1377 *Supplementary Figures*

1378

1379

1380

1381



*AbnU*  
A0A1S1R073

*AbyU*  
F4F7G1

*Cyc15*  
A0A401MXE6

*Pyrl4*  
K7QVW7

1382

1383

1384 **FIGURE S1-** PDB structures of AbnU, AbyU, Cyc15 and Pyrl4

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

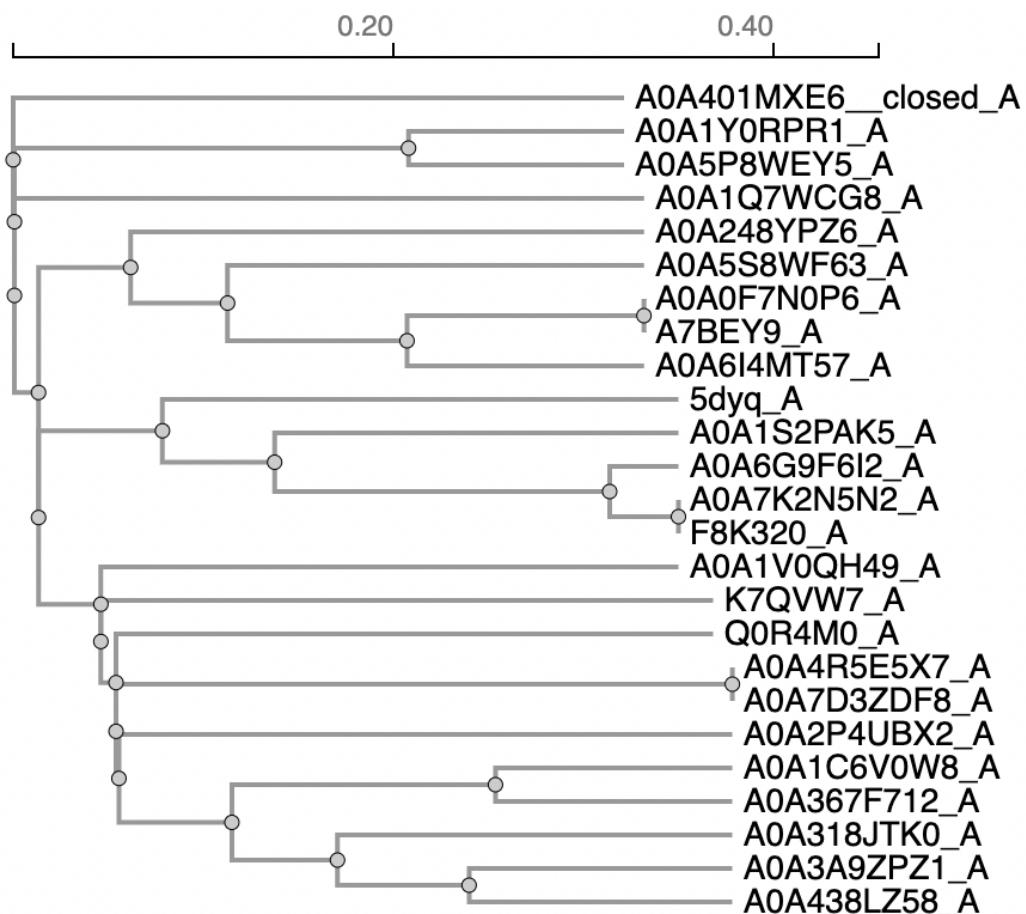
1402

1403

1404

1405

1406



1407

1408 **FIGURE S2-** MSA visualized by Tree (Clustal Omega).

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

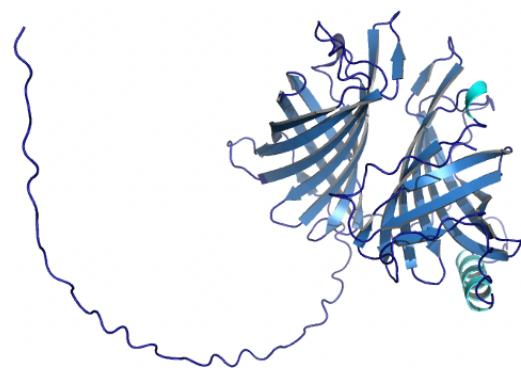
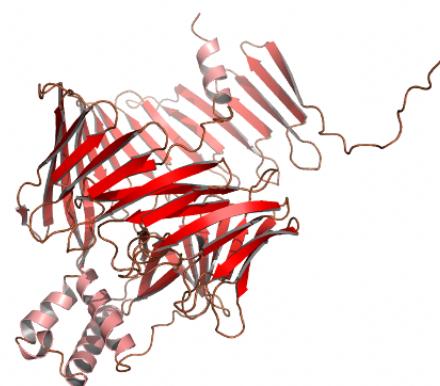
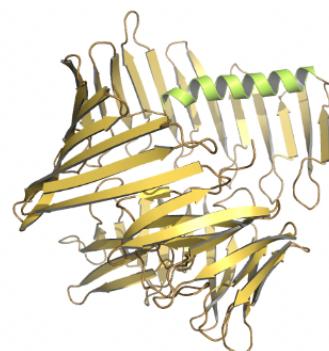
1419

1420

1421

1422

1423  
1424  
1425

**K4FDH3****B**  
**A0A838IM54****C**  
**A0A522WLH4**

1426  
1427  
1428 **FIGURE S3-** Relevant hits that do not display a B-barrel fold  
1429  
1430  
1431

1432

1433