

Wilfrid Laurier University

Scholars Commons @ Laurier

---

Biology Faculty Publications

Biology

---

2022

## CBP60-DB: An AlphaFold-predicted plant kingdom-wide database of the CALMODULIN-BINDING PROTEIN 60 (CBP60) protein family with a novel structural clustering algorithm

Keaun Amani

*Wilfrid Laurier University, aman5230@mylaurier.ca*

Vanessa Shivnauth

*Wilfrid Laurier University, vanessashivnauth@gmail.com*

Christian Castroverde

*dcastroverde@wlu.ca*

Follow this and additional works at: [https://scholars.wlu.ca/biol\\_faculty](https://scholars.wlu.ca/biol_faculty)



Part of the [Bioinformatics Commons](#), [Biology Commons](#), [Computational Biology Commons](#), [Plant Sciences Commons](#), and the [Structural Biology Commons](#)

---

### Recommended Citation

Keaun Amani, Vanessa Shivnauth, Christian Danve M. Castroverde. CBP60-DB: An AlphaFold-predicted plant kingdom-wide database of the CALMODULIN-BINDING PROTEIN 60 (CBP60) protein family with a novel structural clustering algorithm. bioRxiv 2022.07.07.499200; doi: <https://doi.org/10.1101/2022.07.07.499200>

This Article is brought to you for free and open access by the Biology at Scholars Commons @ Laurier. It has been accepted for inclusion in Biology Faculty Publications by an authorized administrator of Scholars Commons @ Laurier. For more information, please contact [scholarscommons@wlu.ca](mailto:scholarscommons@wlu.ca).

1 **CBP60-DB: An AlphaFold-predicted plant kingdom-wide database**  
2 **of the CALMODULIN-BINDING PROTEIN 60 (CBP60) protein**  
3 **family with a novel structural clustering algorithm**

4  
5 Keun Amani\*, Vanessa Shivnauth and Christian Danve M. Castroverde\*

6  
7 **Affiliation:**

8 Department of Biology, Wilfrid Laurier University, Waterloo, ON, Canada N2L 3C5

9  
10 Corresponding Authors: [aman5230@mylaurier.ca](mailto:aman5230@mylaurier.ca); [dcastroverde@wlu.ca](mailto:dcastroverde@wlu.ca)

11

12 **Abstract**

13 Molecular genetic analyses in the model species *Arabidopsis thaliana* have demonstrated the  
14 major roles of different CAM-BINDING PROTEIN 60 (CBP60) proteins in growth, stress  
15 signaling, and immune responses. Prominently, CBP60g and SARD1 are paralogous CBP60  
16 transcription factors that regulate numerous components of the immune system, such as cell  
17 surface and intracellular immune receptors, MAP kinases, WRKY transcription factors, and  
18 biosynthetic enzymes for immunity-activating metabolites salicylic acid (SA) and *N*-  
19 hydroxy-pipecolic acid (NHP). However, their function, regulation and diversification in most  
20 species remain unclear. Here we have created CBP60-DB, a structural and bioinformatic  
21 database that comprehensively characterized 1052 *CBP60* gene homologs (encoding 2376  
22 unique transcripts and 1996 unique proteins) across 62 phylogenetically diverse genomes in the  
23 plant kingdom. We have employed deep learning-predicted structural analyses using AlphaFold2  
24 and then generated dedicated web pages for all plant CBP60 proteins. Importantly, we have  
25 generated a novel clustering visualization algorithm to interrogate kingdom-wide structural  
26 similarities for more efficient inference of conserved functions across various plant taxa.  
27 Because well-characterized CBP60 proteins in *Arabidopsis* are known to be transcription factors  
28 with putative calmodulin-binding domains, we have integrated external bioinformatic resources  
29 to analyze protein domains and motifs. Collectively, we present a plant kingdom-wide  
30 identification of this important protein family in a user-friendly AlphaFold-anchored database,  
31 representing a novel and significant resource for the broader plant biology community.

32

### 33 **Introduction**

34 Plants employ constitutive and inducible defense mechanisms to combat invading pests and  
35 pathogens (Wittstock and Gershenzon, 2002; Freeman, 2008; Zhou and Zhang, 2021). A central  
36 inducible defense response is the production of the plant hormone salicylic acid (SA), which has  
37 essential roles in immunity (Ding & Ding, 2020; Peng et al., 2021; Shields et al., 2022) and  
38 abiotic stress tolerance (Gharbi et al., 2018; Khan et al., 2019; Saleem et al., 2021). Thorough  
39 understanding of plant immunity and stress responses are important in reducing global crop  
40 losses and ensuring food security worldwide (Bailey-Serres et al., 2019; Savary et al., 2019).

41  
42 In the model plant species *Arabidopsis thaliana*, SA production in response to stress is mediated  
43 by the sequential action of the ISOCHORISMATE SYNTHASE 1 (ICS1), ENHANCED  
44 DISEASE SUSCEPTIBILITY 5 (EDS5) and AVRPPHB SUSCEPTIBLE 3 (PBS3) proteins  
45 (Rekhter et al., 2019), which are controlled at the transcriptional level by the master transcription  
46 factor CAM-BINDING PROTEIN 60-LIKE G (CBP60g) and its functionally redundant  
47 homolog SAR Deficient 1 (SARD1; Wang et al., 2009; Zhang et al., 2010; Wang et al., 2011;  
48 Sun et al., 2015). Notably, it is known that SA production and plant immunity are vulnerable to  
49 warming temperatures (Huot et al., 2017; Castroverde and Dina, 2021). This critical temperature-  
50 vulnerability of the plant immune system is controlled via CBP60g/SARD1 (Kim et al., 2022),  
51 which are members of the broadly conserved plant CBP60 protein family (Zheng et al., 2021).

52  
53 Biological understanding of protein function relies on detailed characterization of protein  
54 structures. However, accurate prediction of protein structure from amino acid sequence alone has  
55 remained a central problem in biology (Dill et al., 2008). Traditional methods, such as X-ray  
56 crystallography or NMR spectroscopy, are usually very expensive, time consuming, and can fail  
57 to produce viable results for complexes, membrane-bound proteins, or proteins that are unable to  
58 crystallize (Tugarinov et al., 2004; Shi, 2014; Nogales and Scheres, 2015). A major advance to  
59 solve this grand challenge occurred with the launch of AlphaFold2, which is a novel deep  
60 learning approach for accurately predicting the three-dimensional structure of a protein from its  
61 amino acid sequence (Jumper et al., 2021). However, base AlphaFold2 also suffers from a few  
62 drawbacks, such as lack of exposure for certain internal settings (e.g., number of recycling  
63 steps), it is slightly unoptimized, and the default MSA generation algorithms used can be slow

64 and time-consuming (Mirdita et al., 2022). ColabFold (Mirdita et al., 2022) is an AlphaFold2  
65 derivative that addresses the aforementioned issues with AlphaFold2 and is able to generate  
66 highly accurate predictions comparable, if not superior to those of AlphaFold2. Furthermore,  
67 ColabFold can produce more predictions within a shorter period.

68

69 Because of the biological importance of CBP60g and SARD1 proteins for plant immune system  
70 resilience under changing environmental conditions (Wan et al., 2012; Choudhary et al., 2022;  
71 Kim et al., 2022), it is critical that we fully understand their structures and functions in other  
72 plants. This mechanistic knowledge has critical ramifications on safeguarding plant disease  
73 resistance for a warming climate. Although a recent study conducted a kingdom-wide  
74 phylogenetic analysis of the CBP60 family and potential protein neofunctionalization (Zheng et  
75 al., 2021), there is little functional and molecular information on these proteins in most plants,  
76 including agriculturally important crop species.

77

78 To further understand the diversity of CBP60 protein structure and function in the plant  
79 kingdom, we have created a fully curated, AlphaFold-generated (Jumper et al., 2021) structural  
80 database called the Plant CBP60 Protein Family Database or CBP60-DB  
81 (<https://cbp60db.wlu.ca/>). Of note, this paper describes an algorithm that to our knowledge is a  
82 novel approach to accurately clustering proteins by structural similarity. The proposed algorithm  
83 is simple, accurate, and can be easily reproduced on any modern device. By building our novel  
84 visual clustering algorithm, we were able to compare and cluster the predicted protein structures,  
85 facilitating easier ortholog selection and inference of putative biological functions. A Google  
86 Colaboratory notebook is provided, as well as a minimal implementation for executing locally.  
87 We have showcased a visualization for this algorithm on the index page of the CBP60-DB web  
88 application.

89

### 90 **Plant kingdom-wide sequence collection and AlphaFold-based protein folding**

91 We first identified CBP60 genes and proteins in plant species with published and fully sequenced  
92 genomes. Using the Gramene comparative genomics website (<http://gramene.org/>; Tello-Ruiz et  
93 al., 2020), we obtained a comprehensive kingdom-wide list of representative plant species and  
94 *CBP60* gene homologs in these species. Our base dataset consisted of species names, gene

95 sequences, transcript/cDNA sequences, and protein sequence data. Each protein entry's amino  
96 acid sequence was used as an input to ColabFold for structural predictions. ColabFold  
97 (<https://github.com/sokrypton/ColabFold>) was used instead of the original AlphaFold2 since the  
98 former produces a higher number of predictions within a shorter time, while also improving  
99 prediction quality compared to base AlphaFold2. This improvement is primarily due to  
100 ColabFold's usage of the MMseqs2 algorithm for faster homology search as well as other model  
101 optimizations (Steinegger and Söding, 2017). Furthermore, ColabFold makes some of  
102 AlphaFold2's internal settings easily accessible and configurable, allowing us to adjust settings  
103 such as the number of recycling iterations.

104

### 105 **Database implementation**

106 Prior to the development of CBP60-DB and its web application components, we determined that  
107 an effective solution must be scalable, responsive, simple to use, and sufficiently modular, so  
108 that the application could easily be adapted to other protein families and similar projects. The  
109 final version of our database contains 1996 unique predicted structures (from 2376  
110 corresponding cDNA/transcripts and 1052 unique genes), as well as corresponding metadata, and  
111 confidence metrics. The predicted structures are available in the protein databank (PDB; Berman  
112 et al., 2000) and the newer macromolecular Crystallographic Information File (mmCIF) file  
113 formats.

114

115 The CBP60-DB user interface was designed to be easy to navigate, with an emphasis on several  
116 intuitive visualization options that are available and assembled for best user accessibility.

117 Additionally, the application was written in the Go programming language without third party  
118 dependencies, making it straightforward to re-deploy across any modern system. All database  
119 contents are either stored within the assets directory of the application, which is freely accessible  
120 via HTTP(S), or stored within an internal json file that is then loaded into memory as a hash  
121 table, where keys are the md5 hashes of the unique transcript names. The advantage of an  
122 internal hash table over a traditional database management system (DBMS) is that the internal  
123 hash table is faster for accessing and serving data and requires no additional dependencies.  
124 Furthermore, since the contents of the database are static and the memory required to load the  
125 json file is reasonable (6.9 Mb), there is little need for using an alternative DBMS. However,

126 should we decide to scale the contents of the database to include vastly more entries, an  
127 alternative DBMS will be the preferable solution.

128

### 129 **Data archival**

130 CBP60-DB archives and provides access to the following data below. Note that protein  
131 structures which have been updated, replaced, or removed will not be archived.

- 132 ● Predicted protein crystal structure in PDB and mmCIF file formats.
- 133 ● Protein metadata and AlphaFold2 metadata in json format.
- 134 ● Generated thumbnails of the predicted structure in png format.
- 135 ● AlphaFold2 scoring metrics in json format (pLDDT, PAE, and pTM score).
- 136 ● MMseqs2 MSA file used during model inference in a3m format.
- 137 ● Cluster map of predicted structures in json format.
- 138 ● Phylogenetic tree created within MEGA using the Multiple Sequence Comparison by  
139 Log-Expectation (MUSCLE) alignment algorithm in FASTA format.
- 140 ● Phylogenetic tree generated by FastTree in the Newick file format.

141

142 The predicted Local Distance Difference Test (pLDDT- $C\alpha$ ) is a per residue metric used by  
143 AlphaFold2 to gauge the model's confidence in the position and orientation of each residue  
144 within a predicted structure. Values range from 0 to 100, where higher values are associated with  
145 greater prediction accuracy and less disorder (Jumper et al., 2021).

146

147 The Predicted Aligned Error (PAE) is a  $N_{\text{res}} \times N_{\text{res}}$  matrix where  $N_{\text{res}}$  corresponds to the number  
148 of residues within the input amino acid sequence. Each element within the matrix represents the  
149 predicted distance error in Ångströms of the 1<sup>st</sup> residue's position when aligned on the 2<sup>nd</sup>  
150 residue (Varadi et al., 2021).

151

152 The Molecular Evolutionary Genetics Analysis (Tamura et al., 2021) application was used to  
153 produce the alignment fasta file using the Multiple Sequence Comparison by Log-Expectation  
154 (MUSCLE) (Edgar, 2004) algorithm with the following parameters (Supplementary Table 1).  
155 The alignment fasta file was then used by the Fast Tree algorithm (Price et al., 2009) to produce  
156 a phylogenetic tree in the Newick file format.

157

## 158 **Clustering proteins by structural similarity**

159 Clustering proteins by their structural similarity is an invaluable method for finding proteins with  
160 potentially similar function but with diverging sequences, especially for large protein families  
161 (Mai et al., 2016; Teletin et al., 2019). Traditional sequence-based cluster algorithms also  
162 provide simple and computationally efficient ways of representing similar proteins but have a  
163 major drawback with regards to proteins with similar functionality but different sequences  
164 (Krissinel, 2007; Kosloff and Kolodny, 2007).

165

166 We have proposed a novel algorithm that is simple and effective at clustering proteins by  
167 structural similarity, while also being easily parallelizable. Our algorithm utilizes metrics used  
168 for protein structure comparison (e.g., TM-Align, Root Mean Square Deviation (RMSD), etc.) to  
169 produce a feature tensor that is then used as input to the Uniform Manifold Approximation and  
170 Projection (UMAP; McInnes et al., 2018) algorithm. The corresponding UMAP projection can  
171 then be used as an intuitive visualization, where proteins that are more structurally similar to one  
172 another will be clustered within closer proximity to each other. The advantages of our algorithm  
173 are that it is trivial to implement, easy to utilize, and highly configurable with regards to feature  
174 selection and UMAP hyper-parameter tuning. Furthermore, our algorithm can cluster small  
175 datasets of protein with minimal hardware and within a reasonable amount of time. However, a  
176 drawback of the algorithm is its quadratic time complexity which does not allow it to efficiently  
177 scale on lower-end hardware.

178

179 To produce the input feature tensor pairwise structural comparison, metrics such as TM-Align  
180 (Zhang, 2005) optionally alongside other metrics such as RMSD were used to produce a  
181  $n \times n \times m$  feature tensor, where  $n$  is the number of proteins and  $m$  is the number of features.  
182 The feature tensor was then flattened to produce a  $n \times (n \times m)$  matrix, which was used as an  
183 input to UMAP. UMAP is a powerful dimensionality reduction algorithm that can generally  
184 create more meaningful representations compared to principal component analysis, while also  
185 outperforming t-distributed stochastic neighbor embedding (t-SNE; McInnes et al., 2018). It is  
186 also noteworthy to mention that swapping UMAP with t-SNE produces comparable projections;



187 however, UMAP is significantly faster and, in our opinion, generally produces more intuitive  
188 projections.

189

190 A Google Colaboratory notebook demo

191 (<https://colab.research.google.com/drive/1LOZY33CSO5-PdJAdDApyPlfxUu4DHjcW>) and

192 minimal Python implementation for the clustering algorithm are available. Additionally, a

193 structural cluster of all proteins available within the CBP60-DB is available on the index page of

194 the application (Fig 1).

195



196

197 **Fig 1. Screenshot of the top of the protein structure cluster of the entire CBP60-DB within**  
198 **the Database Visualization section of the index page.**

199

## 200 Data access

201 Information stored on CBP60-DB is hosted by Wilfrid Laurier University, which can be accessed

202 using a modern web browser (<https://cbp60db.wlu.ca/>), through the application programming

203 interface (API), as well as through our public github repository

204 (<https://github.com/KeaunAmani/cbp60db/>). Accessing the database through either means

205 provides access to all data including the predicted structures, metadata, thumbnails, prediction

206 metrics, as well as the cluster map. Additionally, viewing CBP60-DB via your web browser

207 provides access to several interactive and intuitive visualizations, featuring an interactive protein

208 viewer, navigable cluster map (Fig 1), interactive plots for prediction metrics, as well as the top  
209 five most structurally similar proteins (if available).

210

## 211 **Navigating the database**

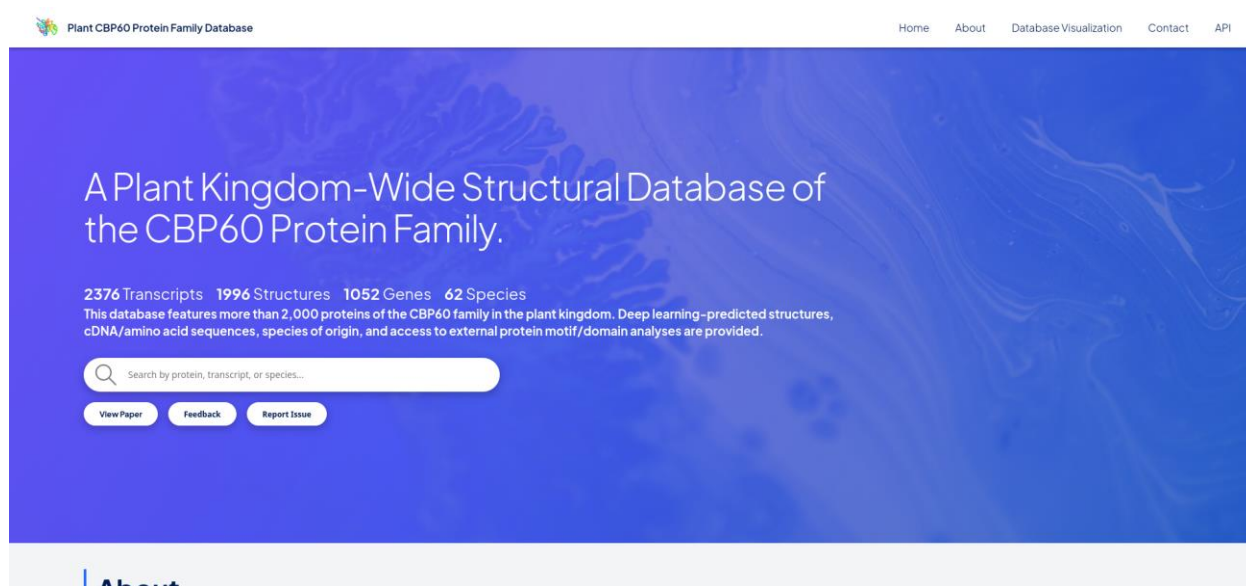
212 There are three primary web pages available on CBP60-DB: (1) index page, (2) protein search  
213 page, and (3) protein information page.

214

### 215 ***Index page:***

216 The CBP60-DB index page (Fig 2) acts as the website home page containing general  
217 information, navigation options, database visualizations, downloads, as well as API endpoint  
218 documentation. To navigate the database, users may either search for a protein directly via the  
219 search bar in the page header, or alternatively interact with the protein structural visualization  
220 cluster. By clicking a node within the cluster, users will be redirected to that protein's  
221 information page. Alternatively, the aforementioned search bar allows users to search for protein  
222 by their transcript name, gene name, or source organism. Proteins that match the search query  
223 will be displayed in the following search page. Another visualization available within this page is  
224 an interactive phylogenetic tree explorer. Note that downloads for the TM-Align Cluster, FASTA  
225 Alignment, and Phylogenetic tree are available underneath their respective visualizations.

226



227

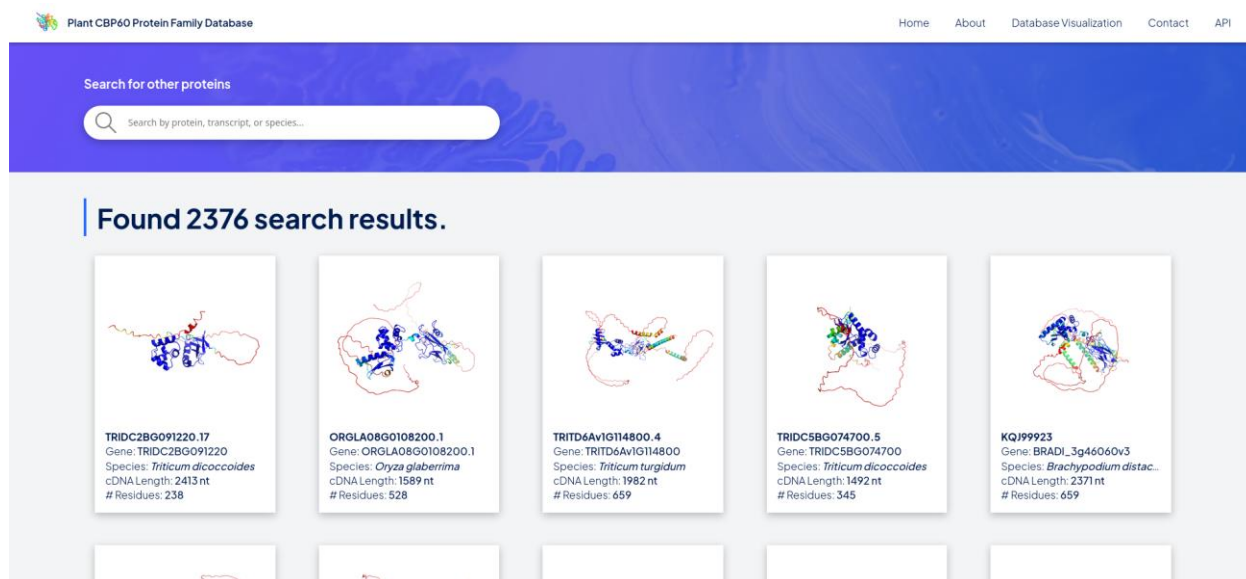
228 **Fig 2. Screenshot of the top of the CBP60-DB index page.**

229

230 ***Search page:***

231 The search page (Fig 3) displays the search results from queries made via the search bar on any  
232 page within CBP60-DB. Once a query is submitted through the search bar, users will be  
233 redirected to this page. If no query is provided, all database entries will be displayed instead.  
234 Search results are unordered and in the form of card previews containing a thumbnail of the  
235 predicted crystal structure, the transcript name, gene name, source organism, cDNA length, and  
236 amino acid sequence length. Users can click on cards to visit their corresponding protein  
237 information pages.

238



239

240 **Fig 3. Screenshot of the top of the CBP60-DB search page.**

241

242 ***Protein information page:***

243 The protein information page (Fig 4) is arguably the most useful page within CBP60-DB,  
244 providing a simple interface for viewing the gene name, transcript name, source organism,  
245 AlphaFold2 settings used, cDNA sequence, amino acid sequence, structure data, redirect to  
246 DNA-Binding Residues tool (Hwang et al., 2007), redirect to Eukaryotic Linear Motifs tool  
247 (Kumar et al., 2020), various downloads as well as visualizations, and the top five most similar  
248 structures (if available) according to the clustering algorithm.

249

250 Data visualizations available on this page include an interactive molecular viewer of the  
251 predicted protein structure utilizing PDBe Molstar (Sehnal et al., 2021), as well as interactive  
252 plots for the PAE and pLDDT scores powered by the plotly.js library (Plotly Technologies Inc.,  
253 2015). Additionally, the exact same protein cluster from the Index page is also available with the  
254 current protein highlighted within the plot. Similar to the Index page, this plot is also navigable  
255 in the same way.

256  
257 Data downloads available on this page consist of the PDB file of the predicted structure, mmCIF  
258 file of the predicted structure, PAE json file, pLDDT json file, amino acid sequence FASTA file,  
259 and the generated MSA used to predict protein structure. These resources are also available for  
260 download directly via the programmatic API.



262  
263 **Fig 4. Screenshot of the top of the CBP60-DB protein information page for the**  
264 **representative protein with the transcript name AT5G26920.1.**

265  
266 **Conclusion and Outlook**

267 Recent in-silico advances for protein structure prediction have accelerated molecular biology  
268 research at an unprecedented scale. Deep learning models have now proven themselves to be  
269 effective tools for protein folding and are made even more valuable through their ease of use and  
270 lower costs compared to traditional techniques (Jumper et al., 2021; Baek et al., 2021). By

271 determining the structures of all proteins within the CBP60 plant kingdom family, biologists can  
272 infer putative functions, evolutionary relationships, and other meaningful information from  
273 protein structures on a broader scale.

274  
275 Overall, the CBP60-DB has generated useful and comprehensive datasets that are foundational  
276 for further functional and molecular studies. Because CBP60 protein family members CBP60  
277 and SARD1 are indispensable master regulators of plant defense responses (Wang et al., 2009;  
278 Zhang et al., 2010; Wang et al., 2011; Sun et al., 2015; Kim et al., 2022), our fundamental  
279 understanding of their structural and functional diversity has profound implications for  
280 mitigating plant diseases. This could potentially address major challenges in agricultural and  
281 natural ecosystems globally, especially on understanding plant immune system resilience  
282 (Velasquez et al., 2018; Kim et al., 2021; Kim et al., 2022) to boost worldwide crop productivity  
283 (Bailey-Serres et al., 2019). Using a robust and rapid bioinformatic pipeline, our comprehensive  
284 deep learning-assisted database with a novel structural clustering algorithm provides the  
285 scientific community with easy-to-access candidate genes/proteins that can be further engineered  
286 to strengthen plant health in a changing world.

287

## 288 **Acknowledgements**

289 Research in the Castroverde Lab is funded by the Natural Sciences and Engineering Research  
290 Council of Canada (NSERC) Discovery Grant (to C.D.M.C.), Canada Foundation for Innovation  
291 (to C.D.M.C.), Ontario Research Fund (to C.D.M.C.), Laurier Faculty of Science institutional  
292 start-up funds (to C.D.M.C.), and a Mitacs Research Training Award (to V.S.). We also thank  
293 Compute Canada and SHARCNET for providing computational power and support, Laurier  
294 Information and Communication Technologies for website hosting, as well as Dr. Sean Johnson  
295 for his critical guidance and feedback.

296

## 297 **References**

298 Ahdritz G, Bouatta N, Kadyan S, Xia Q, Gerecke W, AlQuraishi M. 2021. OpenFold.  
299 doi:[10.5281/zenodo.5709539](https://doi.org/10.5281/zenodo.5709539). <https://github.com/aqlaboratory/openfold>.

300 Baek M, DiMaio F, Anishchenko I, Dauparas J, Ovchinnikov S, Lee GR, Wang J, Cong Q,  
301 Kinch LN, Schaeffer RD, et al. 2021. Accurate prediction of protein structures and interactions  
302 using a three-track neural network. *Science*. 373(6557):871–876. doi:[10.1126/science.abj8754](https://doi.org/10.1126/science.abj8754).

303 Bailey-Serres J, Parker JE, Ainsworth EA, Oldroyd GED, Schroeder JI. 2019. Genetic strategies  
304 for improving crop yields. *Nature*. 575(7781):109–118. doi:[10.1038/s41586-019-1679-0](https://doi.org/10.1038/s41586-019-1679-0).  
305 <https://www.nature.com/articles/s41586-019-1679-0>.

306 Berman HM. 2000. The Protein Data Bank. *Nucleic Acids Research*. 28(1):235–242.  
307 doi:[10.1093/nar/28.1.235](https://doi.org/10.1093/nar/28.1.235).

308 Castroverde CDM, Dina D. 2021. Temperature regulation of plant hormone signaling during  
309 stress and development. *Journal of Experimental Botany*. doi:[10.1093/jxb/erab257](https://doi.org/10.1093/jxb/erab257).

310 Cheng S, Wu R, Yu Z, Li B, Zhang X, Peng J, You Y. 2022. FastFold: Reducing AlphaFold  
311 Training Time from 11 Days to 67 Hours. arXiv:220300854 [cs, q-bio].  
312 <https://arxiv.org/abs/2203.00854>.

313 Choudhary A, Senthil-Kumar M. 2022. Drought attenuates plant defence against bacterial  
314 pathogens by suppressing the expression of *CBP60g* / *SARD1* during combined stress. *Plant, Cell*  
315 *& Environment*. 45(4):1127–1145. doi:[10.1111/pce.14275](https://doi.org/10.1111/pce.14275).

316 Dill KA, Ozkan SB, Shell MS, Weikl TR. 2008. The Protein Folding Problem. *Annual Review*  
317 *of Biophysics*. 37(1):289–316. doi:[10.1146/annurev.biophys.37.092707.153558](https://doi.org/10.1146/annurev.biophys.37.092707.153558).

318 Ding P, Ding Y. 2020 Feb. Stories of Salicylic Acid: A Plant Defense Hormone. *Trends in Plant*  
319 *Science*. doi:[10.1016/j.tplants.2020.01.004](https://doi.org/10.1016/j.tplants.2020.01.004).

320 Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high  
321 throughput. *Nucleic Acids Research*. 32(5):1792–1797. doi:[10.1093/nar/gkh340](https://doi.org/10.1093/nar/gkh340).

322 Freeman. 2008. An Overview of Plant Defenses against Pathogens and Herbivores. *The Plant*  
323 *Health Instructor*. doi:[10.1094/phi-i-2008-0226-01](https://doi.org/10.1094/phi-i-2008-0226-01).

324 Gharbi E, Lutts S, Dailly H, Quinet M. 2018. Comparison between the impacts of two different  
325 modes of salicylic acid application on tomato (*Solanum lycopersicum*) responses to salinity.  
326 *Plant Signaling & Behavior*. 13(5):e1469361. doi:[10.1080/15592324.2018.1469361](https://doi.org/10.1080/15592324.2018.1469361).

327 Huot B, Castroverde CDM, Velásquez AC, Hubbard E, Pulman JA, Yao J, Childs KL, Tsuda K,  
328 Montgomery BL, He SY. 2017. Dual impact of elevated temperature on plant defence and  
329 bacterial virulence in *Arabidopsis*. *Nature Communications*. 8(1). doi:[10.1038/s41467-017-](https://doi.org/10.1038/s41467-017-01674-2)  
330 [01674-2](https://doi.org/10.1038/s41467-017-01674-2).

331 Hwang S, Gou Z, Kuznetsov IB. 2007. DP-Bind: a web server for sequence-based prediction of  
332 DNA-binding residues in DNA-binding proteins. *Bioinformatics*. 23(5):634–636.  
333 doi:[10.1093/bioinformatics/btl672](https://doi.org/10.1093/bioinformatics/btl672).

334 Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O, Tunyasuvunakool K, Bates  
335 R, Židek A, Potapenko A, et al. 2021. Highly accurate protein structure prediction with  
336 AlphaFold. *Nature*. 596(7873):583–589. doi:[10.1038/s41586-021-03819-2](https://doi.org/10.1038/s41586-021-03819-2).

337 Khan A, Kamran M, Imran M, Al-Harrasi A, Al-Rawahi A, Al-Amri I, Lee I-J, Khan AL. 2019.  
338 Silicon and salicylic acid confer high-pH stress tolerance in tomato seedlings. *Scientific Reports*.  
339 9(1). doi:[10.1038/s41598-019-55651-4](https://doi.org/10.1038/s41598-019-55651-4).

340 Kim JH, Castroverde CDM, Huang S, Li C, Hilleary R, Seroka A, Sohrabi R, Medina-Yerena D,  
341 Huot B, Wang J, et al. 2022 Jun 29. Increasing the resilience of plant immunity to a warming  
342 climate. *Nature*. doi:[10.1038/s41586-022-04902-y](https://doi.org/10.1038/s41586-022-04902-y).

343 Kim JH, Hilleary R, Seroka A, He SY. 2021. Crops of the future: building a climate-resilient  
344 plant immune system. *Current Opinion in Plant Biology*. 60:101997.  
345 doi:[10.1016/j.pbi.2020.101997](https://doi.org/10.1016/j.pbi.2020.101997).

346 Kosloff M, Kolodny R. 2007. Sequence-similar, structure-dissimilar protein pairs in the PDB.  
347 *Proteins: Structure, Function, and Bioinformatics*. 71(2):891–902. doi:[10.1002/prot.21770](https://doi.org/10.1002/prot.21770).

348 Krissinel E. 2007. On the relationship between sequence and structure similarities in proteomics.  
349 *Bioinformatics*. 23:717–723. doi:[10.1093/bioinformatics/btm006](https://doi.org/10.1093/bioinformatics/btm006).



350 Kumar M, Gouw M, Michael S, Sámano-Sánchez H, Pancsa R, Glavina J, Diakogianni A,  
351 Valverde JA, Bukirova D, Čalyševa J, et al. 2020. ELM-the eukaryotic linear motif resource in  
352 2020. *Nucleic Acids Research*. 48(D1):D296–D306. doi:[10.1093/nar/gkz1030](https://doi.org/10.1093/nar/gkz1030).

353 Mai T-L, Hu G-M, Chen C-M. 2016. Visualizing and Clustering Protein Similarity Networks:  
354 Sequences, Structures, and Functions. *Journal of Proteome Research*. 15(7):2123–2131.  
355 doi:[10.1021/acs.jproteome.5b01031](https://doi.org/10.1021/acs.jproteome.5b01031).

356 McInnes L, Healy J, Melville J. 2018. UMAP: Uniform Manifold Approximation and Projection  
357 for Dimension Reduction. arXivorg. <https://arxiv.org/abs/1802.03426>.

358 Mirdita M, Schütze K, Moriwaki Y, Heo L, Ovchinnikov S, Steinegger M. 2022. ColabFold:  
359 making protein folding accessible to all. *Nature Methods*.:1–4. doi:[10.1038/s41592-022-01488-](https://doi.org/10.1038/s41592-022-01488-1)  
360 [1](https://doi.org/10.1038/s41592-022-01488-1).

361 Nogales E, Scheres Sjors HW. 2015. Cryo-EM: A Unique Tool for the Visualization of  
362 Macromolecular Complexity. *Molecular Cell*. 58(4):677–689. doi:[10.1016/j.molcel.2015.02.019](https://doi.org/10.1016/j.molcel.2015.02.019).

363 Peng Y, Yang J, Li X, Zhang Y. 2021. Salicylic Acid: Biosynthesis and Signaling. *Annual*  
364 *Review of Plant Biology*. 72:761–791. doi:[10.1146/annurev-arplant-081320-092855](https://doi.org/10.1146/annurev-arplant-081320-092855).

365 Plotly Technologies Inc. 2015. Collaborative data science. <https://plot.ly>.

366 Price MN, Dehal PS, Arkin AP. 2009. FastTree: Computing Large Minimum Evolution Trees  
367 with Profiles instead of a Distance Matrix. *Molecular Biology and Evolution*. 26(7):1641–1650.  
368 doi:[10.1093/molbev/msp077](https://doi.org/10.1093/molbev/msp077).

369 Rekhter D, Lüdke D, Ding Y, Feussner K, Zienkiewicz K, Lipka V, Wiermer M, Zhang Y,  
370 Feussner I. 2019. Isochorismate-derived biosynthesis of the plant stress hormone salicylic acid.  
371 *Science*. 365(6452):498–502. doi:[10.1126/science.aaw1720](https://doi.org/10.1126/science.aaw1720).

372 Saleem M, Fariduddin Q, Castroverde CDM. 2021. Salicylic acid: A key regulator of redox  
373 signalling and plant immunity. *Plant Physiology and Biochemistry*. 168:381–397.  
374 doi:[10.1016/j.plaphy.2021.10.011](https://doi.org/10.1016/j.plaphy.2021.10.011).



375 Savary S, Willocquet L, Pethybridge SJ, Esker P, McRoberts N, Nelson A. 2019. The global  
376 burden of pathogens and pests on major food crops. *Nature Ecology & Evolution*. 3(3):430–439.  
377 doi:[10.1038/s41559-018-0793-y](https://doi.org/10.1038/s41559-018-0793-y).

378 Sehnal D, Bittrich S, Deshpande M, Svobodová R, Berka K, Bazgier V, Velankar S, Burley SK,  
379 Koča J, Rose AS. 2021. Mol\* Viewer: modern web app for 3D visualization and analysis of  
380 large biomolecular structures. GitHub. <https://github.com/molstar/molstar>.

381 Shi Y. 2014. A Glimpse of Structural Biology through X-Ray Crystallography. *Cell*. 159(5):995–  
382 1014. doi:[10.1016/j.cell.2014.10.051](https://doi.org/10.1016/j.cell.2014.10.051).

383 Shields A, Shivnauth V, Castroverde CDM. 2022. Salicylic Acid and N-Hydroxypipicolinic Acid  
384 at the Fulcrum of the Plant Immunity-Growth Equilibrium. *Frontiers in Plant Science*. 13.  
385 doi:[10.3389/fpls.2022.841688](https://doi.org/10.3389/fpls.2022.841688).

386 Steinegger M, Söding J. 2017. MMseqs2 enables sensitive protein sequence searching for the  
387 analysis of massive data sets. *Nature Biotechnology*. 35(11):1026–1028. doi:[10.1038/nbt.3988](https://doi.org/10.1038/nbt.3988).

388 Sun T, Zhang Y, Li Y, Zhang Q, Ding Y, Zhang Y. 2015. ChIP-seq reveals broad roles of  
389 SARD1 and CBP60g in regulating plant immunity. *Nature Communications*. 6(1).  
390 doi:[10.1038/ncomms10159](https://doi.org/10.1038/ncomms10159).

391 Tamura K, Stecher G, Kumar S. 2021. MEGA11: Molecular evolutionary genetics analysis  
392 version 11. *Molecular Biology and Evolution*. 38:3022–3027. doi:[10.1093/molbev/msab120](https://doi.org/10.1093/molbev/msab120).

393 Teletin M, Czibula G, Bocicor M-I. 2019. Using clustering models for uncovering proteins’  
394 structural similarity. *IEEE Xplore*.:185–190. doi:[10.1109/SACI46893.2019.9111642](https://doi.org/10.1109/SACI46893.2019.9111642).  
395 <https://ieeexplore.ieee.org/document/9111642>.

396 Tello-Ruiz MK, Naithani S, Gupta P, Olson A, Wei S, Preece J, Jiao Y, Wang B, Chougule K,  
397 Garg P, et al. 2020. Gramene 2021: harnessing the power of comparative genomics and  
398 pathways for plant research. *Nucleic Acids Research*. 49(D1):D1452–D1463.  
399 doi:[10.1093/nar/gkaa979](https://doi.org/10.1093/nar/gkaa979).

400 Tugarinov V, Hwang PM, Kay LE. 2004. Nuclear Magnetic Resonance Spectroscopy of High-  
401 Molecular-Weight Proteins. *Annual Review of Biochemistry*. 73(1):107–146.  
402 doi:[10.1146/annurev.biochem.73.011303.074004](https://doi.org/10.1146/annurev.biochem.73.011303.074004).

403 Varadi M, Anyango S, Deshpande M, Nair S, Natassia C, Yordanova G, Yuan D, Stroe O, Wood  
404 G, Laydon A, et al. 2021. AlphaFold Protein Structure Database: massively expanding the  
405 structural coverage of protein-sequence space with high-accuracy models. *Nucleic Acids*  
406 *Research*. doi:[10.1093/nar/gkab1061](https://doi.org/10.1093/nar/gkab1061).

407 Wan D, Li R, Zou B, Zhang X, Cong J, Wang R, Xia Y, Li G. 2012. Calmodulin-binding protein  
408 CBP60g is a positive regulator of both disease resistance and drought tolerance in Arabidopsis.  
409 *Plant Cell Reports*. 31(7):1269–1281. doi:[10.1007/s00299-012-1247-7](https://doi.org/10.1007/s00299-012-1247-7).

410 Wang L, Tsuda K, Sato M, Cohen JD, Katagiri F, Glazebrook J. 2009. Arabidopsis CaM Binding  
411 Protein CBP60g Contributes to MAMP-Induced SA Accumulation and Is Involved in Disease  
412 Resistance against *Pseudomonas syringae*. *PLoS Pathogens*. 5(2):e1000301.  
413 doi:[10.1371/journal.ppat.1000301](https://doi.org/10.1371/journal.ppat.1000301).

414 Wang L, Tsuda K, Truman W, Sato M, Nguyen LV, Katagiri F, Glazebrook J. 2011. CBP60g  
415 and SARD1 play partially redundant critical roles in salicylic acid signaling. *The Plant Journal*.  
416 67(6):1029–1041. doi:[10.1111/j.1365-313x.2011.04655.x](https://doi.org/10.1111/j.1365-313x.2011.04655.x).

417 Wittstock U, Gershenzon J. 2002. Constitutive plant toxins and their role in defense against  
418 herbivores and pathogens. *Current Opinion in Plant Biology*. 5(4):300–307. doi:[10.1016/s1369-  
419 5266\(02\)00264-9](https://doi.org/10.1016/s1369-5266(02)00264-9).

420 Zhang Y. 2005. TM-align: a protein structure alignment algorithm based on the TM-score.  
421 *Nucleic Acids Research*. 33(7):2302–2309. doi:[10.1093/nar/gki524](https://doi.org/10.1093/nar/gki524).

422 Zhang Y, Xu S, Ding P, Wang D, Cheng YT, He J, Gao M, Xu F, Li Y, Zhu Z, et al. 2010.  
423 Control of salicylic acid synthesis and systemic acquired resistance by two members of a plant-  
424 specific family of transcription factors. *Proceedings of the National Academy of Sciences*.  
425 107(42):18220–18225. doi:[10.1073/pnas.1005225107](https://doi.org/10.1073/pnas.1005225107).

426 Zheng Q, Majsec K, Katagiri F. 2021 Oct 5. Pathogen-driven coevolution across the CBP60  
427 plant immune regulator subfamilies confers resilience on the regulator module. *New Phytologist*.  
428 doi:[10.1111/nph.17769](https://doi.org/10.1111/nph.17769).

429 Zhou J-M, Zhang Y. 2020. Plant Immunity: Danger Perception and Signaling. *Cell*. 181(5):978–  
430 989. doi:[10.1016/j.cell.2020.04.028](https://doi.org/10.1016/j.cell.2020.04.028).

431  
432 **Supplementary Table 1:** Configuration used for MUSCLE alignment generation using  
433 MEGA11.

434

Option	Value used
Gap Open Penalty	-2.90
Gap Extend Penalty	0.00
Hydrophobicity Multiplier Penalty	1.20
Max Memory in MB	4096
Max Iterations	16
Cluster Method (Iterations 1,2)	UPGMA
Cluster Method (Other Iterations)	UPGMA
Min Diag Length (Lambda)	24

435