

DE-STRESS: A user-friendly web application for the evaluation of protein designs

Michael J. Stam¹ and Christopher W. Wood^{2*}.

1. School of Informatics, University of Edinburgh, 10 Crichton St, Newington, Edinburgh EH8 9AB.

2. School of Biological Sciences, University of Edinburgh, Roger Land Building, Edinburgh EH9 3FF.

* To whom correspondence should be addressed.

Abstract

Motivation

It is becoming routine to design protein structures *de novo*, with many interesting and useful examples in the literature. However, most sequences of designed proteins could be classed as failures when characterised in the lab, usually as a result of low expression, misfolding, aggregation or lack of function. This high attrition rate makes protein design unreliable and costly. These limitations could potentially be addressed if it were quick and easy to generate a set of high-quality metrics and information regarding designs, which could be used to make reproducible and data-driven decisions about which designs to characterise experimentally.

Results

We present DE-STRESS (DEsigned STRucture Evaluation ServiceS), a web application for the evaluation of structural models of designed and engineered proteins. DE-STRESS has been designed to be simple, intuitive to use and responsive. It provides a wealth of information regarding designs, as well as tools to help contextualise the results and formally describe the properties that a design requires to be fit for purpose.

Availability

DE-STRESS is available for non-commercial use, without registration, through the following website: <https://pragmaticproteindesign.bio.ed.ac.uk/de-stress/>. Source code for the

28 application is available on GitHub: <https://github.com/wells-wood-research/de-stress>. The
29 data used to generate reference sets is available through a GraphQL API, with the following
30 URL: <https://pragmaticproteindesign.bio.ed.ac.uk/big-structure/graphql>.

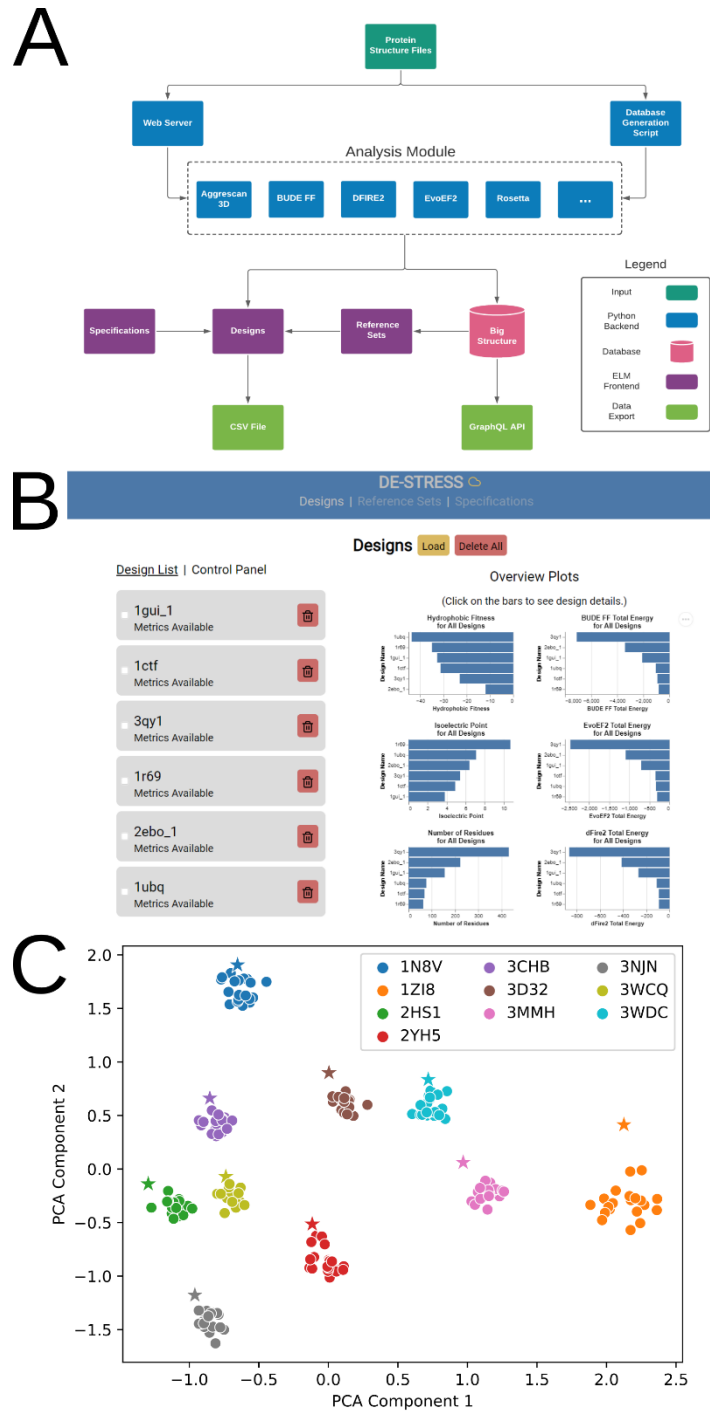
31 Introduction

32 There has been rapid development in the field of *de novo* protein design over recent years,
33 with more groups producing increasingly ambitious designs with complex behaviour, often
34 applied in cellular environments (Ben-Sasson et al., 2021; Glasgow et al., 2019; Harrington
35 et al., 2021; Herud-Sikimić et al., 2021; Pirro et al., 2020; Sesterhenn et al., 2020; VanDrise
36 et al., 2021).

37 Despite the great promise of *de novo* protein design, it remains the domain of highly
38 specialist research groups, as there are significant barriers blocking broader adoption as a
39 methodology. One major challenge is that only a fraction of designs adopt stable, folded
40 structures when expressed (Huang et al., 2016), and it can be difficult to identify these
41 models using the metrics calculated during the design process alone (Radom et al., 2018).
42 This is especially challenging for designs with complex requirements that are needed for
43 targeted applications.

44 Here we present DE-STRESS (DEsigned STRucture Evaluation ServiceS), a user-friendly
45 web application for evaluating structural models of designed and engineered proteins. We
46 aim to provide the user with as much information as possible about their designs before they
47 select sequences to characterise experimentally.

48 Methods and Results



49 *Figure 1: Overview of the DE-STRESS application. A) Architecture of the application. B) The*
 50 *“Designs” page. C) Principal component analysis of DE-STRESS metrics generated for*
 51 *experimentally-determined structures (stars) and folding decoys (circles).*

52 The DE-STRESS application consists of a simple and intuitive user interface, written in
 53 Elm/JavaScript, and a backend web stack, consisting of

54 Unicorn/Flask/GraphQL/PostgreSQL (figure 1A). The interface has three main sections that
55 the user can explore: Designs, Reference Sets and Specifications.

56 On the Designs page, users can upload models of proteins (in PDB format) to the DE-
57 STRESS server, where all the included metrics will be calculated for each design. Once the
58 metrics have been calculated, an overview of the whole batch of designs is provided on the
59 front page (figure 1B). Detailed information can be viewed for each design, as well as a
60 comparison to the active Reference Set and Requirement Specification (*vide infra*).

61 On the Reference Sets page, users can define a set of known protein structures from the
62 PDB (Berman et al., 2003), which can be used as a basis of comparison for their designs.
63 We have precalculated the metrics included in DE-STRESS for the biological units of 82,010
64 protein structures, as defined by the PDBe

65 (<http://ftp.ebi.ac.uk/pub/databases/pdb/data/biounit/>). The remaining structures in the PDB
66 either did not contain protein, contained formatting errors in the PDB file or, in the case of
67 large structures, failed to return results within a reasonable timeframe. Using these data, the
68 user can define their own reference sets by submitting a list of PDB accession codes,
69 enabling them to compare their designs to relevant structures. Additionally, two default
70 reference sets are provided as an example, based on high-quality structures from Top500
71 (Hobohm and Sander, 1994) and Pisces (Wang and Dunbrack, 2003). Once a reference set
72 has been defined, aggregated metrics are presented alongside the metrics for the user's
73 designs. All the data used to generate the reference sets is available to search and
74 download, programmatically and interactively, through a GraphQL API available at the
75 following url: <https://pragmaticproteindesign.bio.ed.ac.uk/big-structure/graphql>.

76 Finally, the Specifications page allows the user to define "Requirement Specifications",
77 which encapsulate the properties their designs should have in order to be fit for purpose.
78 The user can define complex rules that can be used to filter designs, alongside associated
79 metadata. We plan to expand the role of the specifications in the future, allowing the user to

80 capture more information about their design intent and export the specification to be used by
81 other programmes.

82 A variety of external software packages are used by the DE-STRESS web server to
83 calculate metrics for uploaded protein structures. Basic information about the protein
84 structure is extracted using ISAMBARD (Wood et al., 2017), including information such as
85 the isoelectric point and composition of the sequence, as well as implementations of a few
86 metrics from the literature, such as packing density (Weiss, 2007) and hydrophobic fitness
87 (Huang et al., 1995). In addition to these metrics, DE-STRESS applies a range of scoring
88 functions that are well established in the protein-design field, such as BUDE (McIntosh-
89 Smith et al., 2012, 2015), EvoEF2 (Huang et al., 2020), Rosetta (Alford et al., 2017) and
90 DFIRE2 (Yang and Zhou, 2008). Finally, Aggrescan3D (Kuriata et al., 2019) calculates an
91 aggregation propensity score for protein structures. Additional metrics will be incorporated
92 into DE-STRESS in future releases. These metrics are presented on the Design Details
93 page, alongside a visualisation of the model, using the NGL JavaScript library (Rose and
94 Hildebrand, 2015; Rose et al., 2016), and other information such as secondary structure
95 assignment using DSSP (Kabsch and Sander, 1983; Touw et al., 2015).

96 A privacy first approach has been taken when implementing DE-STRESS. No login is
97 required to use the application and no data regarding the user, or their designs, are stored
98 on our server. Designs are submitted directly to an in-memory job queue, with no associated
99 metadata, and the results are returned directly to the user. All data regarding the user's
100 designs are stored locally on the device used to access the website and can be exported to
101 a CSV file for further analysis. With this architecture, we aim to give the user confidence in
102 submitting their designs to the server. However, if they would like to take further steps to
103 ensure that no one could access their data, they can run a local instance of the web
104 application, which we have made as simple as possible by containerising the application.

105 We envisage that DE-STRESS will be useful for generating descriptive information and
106 statistics that could be manually examined by users to choose designs that meet the needs

107 of their application. Beyond this, the datasets that DE-STRESS creates could be useful for
108 automatic identification of high-quality designs using data-driven methods. As a simple
109 example of this, we attempted to identify folding decoys from experimentally-determined
110 structures. Using the DE-STRESS web application, we generated and exported metrics for a
111 random sample of 10 experimentally-determined structures, along with 200 decoys (20 per
112 structure) generated by 3DRobot (Deng et al., 2016). The metrics were normalised, using
113 min-max scaling, and principal component analysis was performed. After this, the first two
114 principal components were plotted against each other, and the experimentally-determined
115 structure, with their associated decoys, formed neat clusters (figure 1C). Furthermore, the
116 experimentally-determined structures were close to, but distinct from, the main cluster,
117 indicating that the metrics included in DE-STRESS could be used to automatically identify
118 high-quality models using machine learning. The dataset and the associated scripts for
119 performing this analysis are available on GitHub: [https://github.com/wells-wood-](https://github.com/wells-wood-research/stam-m-wood-c-de-stress-2021)
120 [research/stam-m-wood-c-de-stress-2021](https://github.com/wells-wood-research/stam-m-wood-c-de-stress-2021).

121 Conclusions

122 DE-STRESS enables both non-experts and seasoned protein designers to rapidly evaluate
123 their designs, providing a framework for making reproducible, data-driven decisions about
124 which design to take forward for experimental characterisation. While some protocols and
125 applications have been developed to address some of the same challenges as DE-STRESS
126 (Bernhofer et al., 2021; Guffy et al., 2018; Yallapragada et al., 2020), none of them have the
127 same breadth of metrics and tools, all packaged in a user-friendly web application.

128 It is our aim that using DE-STRESS will reduce the failure rate of designs taken into the lab,
129 thus increasing the efficiency of protein design, making it more accessible and reliable as a
130 technique.

131 Acknowledgements

132 The authors would like to thank Lynne Regan and Dek Woolfson for feedback on the
133 manuscript and application, the UoE School of Biological Sciences IT department for
134 infrastructure support, and the developers of the software used to generate many of the
135 metrics included in DE-STRESS.

136 Funding

137 CWW is supported by an Engineering and Physical Sciences Research Council Fellowship
138 (EP/S003002/1). Michael Stam is supported by the United Kingdom Research and
139 Innovation (grant EP/S02431X/1), UKRI Centre for Doctoral Training in Biomedical AI at the
140 University of Edinburgh, School of Informatics. This work was supported by the Wellcome
141 Trust-University of Edinburgh Institutional Strategic Support Fund (ISSF3).

142 References

- 143 Alford, R.F., Leaver-Fay, A., Jeliaskov, J.R., O'Meara, M.J., DiMaio, F.P., Park, H.,
144 Shapovalov, M.V., Renfrew, P.D., Mulligan, V.K., Kappel, K., et al. (2017). The Rosetta All-
145 Atom Energy Function for Macromolecular Modeling and Design. *J. Chem. Theory Comput.*
146 *13*, 3031–3048.
- 147 Ben-Sasson, A.J., Watson, J.L., Sheffler, W., Johnson, M.C., Bittleston, A., Somasundaram,
148 L., Decarreau, J., Jiao, F., Chen, J., Mela, I., et al. (2021). Design of biologically active
149 binary protein 2D materials. *Nature* *589*, 468–473.
- 150 Berman, H., Henrick, K., and Nakamura, H. (2003). Announcing the worldwide Protein Data
151 Bank. *Nature Structural & Molecular Biology* *10*, 980–980.
- 152 Bernhofer, M., Dallago, C., Karl, T., Satagopam, V., Heinzinger, M., Littmann, M., Olenyi, T.,
153 Qiu, J., Schütze, K., Yachdav, G., et al. (2021). PredictProtein – Predicting Protein Structure
154 and Function for 29 Years. *BioRxiv* 2021.02.23.432527.
- 155 Deng, H., Jia, Y., and Zhang, Y. (2016). 3DRobot: automated generation of diverse and well-
156 packed protein structure decoys. *Bioinformatics* *32*, 378–387.
- 157 Glasgow, A.A., Huang, Y.-M., Mandell, D.J., Thompson, M., Ritterson, R., Loshbaugh, A.L.,
158 Pellegrino, J., Krivacic, C., Pache, R.A., Barlow, K.A., et al. (2019). Computational design of
159 a modular protein sense-response system. *Science* *366*, 1024–1028.
- 160 Guffy, S.L., Teets, F.D., Langlois, M.I., and Kuhlman, B. (2018). Protocols for Requirement-
161 Driven Protein Design in the Rosetta Modeling Program. *J. Chem. Inf. Model.* *58*, 895–901.

- 162 Harrington, L., Fletcher, J.M., Heermann, T., Woolfson, D.N., and Schwille, P. (2021). De
163 novo design of a reversible phosphorylation-dependent switch for membrane targeting.
164 *Nature Communications* 12, 1472.
- 165 Herud-Sikimić, O., Stiel, A.C., Kolb, M., Shanmugaratnam, S., Berendzen, K.W., Feldhaus,
166 C., Höcker, B., and Jürgens, G. (2021). A biosensor for the direct visualization of auxin.
167 *Nature* 1–5.
- 168 Hobohm, U., and Sander, C. (1994). Enlarged representative set of protein structures.
169 *Protein Science* 3, 522–524.
- 170 Huang, E.S., Subbiah, S., and Levitt, M. (1995). Recognizing native folds by the
171 arrangement of hydrophobic and polar residues. *J Mol Biol* 252, 709–720.
- 172 Huang, P.-S., Boyken, S.E., and Baker, D. (2016). The coming of age of de novo protein
173 design. *Nature* 537, 320–327.
- 174 Huang, X., Pearce, R., and Zhang, Y. (2020). EvoEF2: accurate and fast energy function for
175 computational protein design. *Bioinformatics* 36, 1135–1142.
- 176 Kabsch, W., and Sander, C. (1983). Dictionary of protein secondary structure: Pattern
177 recognition of hydrogen-bonded and geometrical features. *Biopolymers* 22, 2577–2637.
- 178 Kuriata, A., Iglesias, V., Kurcinski, M., Ventura, S., and Kmiecik, S. (2019). Aggrescan3D
179 standalone package for structure-based prediction of protein aggregation properties.
180 *Bioinformatics* 35, 3834–3835.
- 181 McIntosh-Smith, S., Wilson, T., Ibarra, A.Á., Crisp, J., and Sessions, R.B. (2012).
182 Benchmarking Energy Efficiency, Power Costs and Carbon Emissions on Heterogeneous
183 Systems. *The Computer Journal* 55, 192–205.
- 184 McIntosh-Smith, S., Price, J., Sessions, R.B., and Ibarra, A.A. (2015). High performance in
185 silico virtual drug screening on many-core processors. *The International Journal of High*
186 *Performance Computing Applications* 29, 119–134.
- 187 Pirro, F., Schmidt, N., Lincoff, J., Widel, Z.X., Polizzi, N.F., Liu, L., Therien, M.J., Grabe, M.,
188 Chino, M., Lombardi, A., et al. (2020). Allosteric cooperation in a de novo-designed two-
189 domain protein. *PNAS* 117, 33246–33253.
- 190 Radom, F., Plückthun, A., and Paci, E. (2018). Assessment of ab initio models of protein
191 complexes by molecular dynamics. *PLOS Computational Biology* 14, e1006182.
- 192 Rose, A.S., and Hildebrand, P.W. (2015). NGL Viewer: a web application for molecular
193 visualization. *Nucleic Acids Research* 43, W576–W579.
- 194 Rose, A.S., Bradley, A.R., Valasatava, Y., Duarte, J.M., Prlić, A., and Rose, P.W. (2016).
195 Web-based molecular graphics for large complexes. In *Proceedings of the 21st International*
196 *Conference on Web3D Technology*, (New York, NY, USA: Association for Computing
197 Machinery), pp. 185–186.
- 198 Sesterhenn, F., Yang, C., Bonet, J., Cramer, J.T., Wen, X., Wang, Y., Chiang, C.-I., Abriata,
199 L.A., Kucharska, I., Castoro, G., et al. (2020). De novo protein design enables the precise
200 induction of RSV-neutralizing antibodies. *Science* 368.

- 201 Touw, W.G., Baakman, C., Black, J., te Beek, T.A.H., Krieger, E., Joosten, R.P., and Vriend,
202 G. (2015). A series of PDB-related databanks for everyday needs. *Nucleic Acids Research*
203 43, D364–D368.
- 204 VanDrissse, C.M., Lipsh-Sokolik, R., Khersonsky, O., Fleishman, S.J., and Newman, D.K.
205 (2021). Computationally designed pyocyanin demethylase acts synergistically with
206 tobramycin to kill recalcitrant *Pseudomonas aeruginosa* biofilms. *PNAS* 118.
- 207 Wang, G., and Dunbrack, R.L., Jr (2003). PISCES: a protein sequence culling server.
208 *Bioinformatics* 19, 1589–1591.
- 209 Weiss, M.S. (2007). On the interrelationship between atomic displacement parameters
210 (ADPs) and coordinates in protein structures. *Acta Crystallogr D Biol Crystallogr* 63, 1235–
211 1242.
- 212 Wood, C.W., Heal, J.W., Thomson, A.R., Bartlett, G.J., Ibarra, A.Á., Brady, R.L., Sessions,
213 R.B., and Woolfson, D.N. (2017). ISAMBARD: an open-source computational environment
214 for biomolecular analysis, modelling and design. *Bioinformatics* 33, 3043–3050.
- 215 Yallapragada, V.V.B., Walker, S.P., Devoy, C., Buckley, S., Flores, Y., and Tangney, M.
216 (2020). Function2Form Bridge—Toward synthetic protein holistic performance prediction.
217 *Proteins: Structure, Function, and Bioinformatics* 88, 462–475.
- 218 Yang, Y., and Zhou, Y. (2008). Ab initio folding of terminal segments with secondary
219 structures reveals the fine difference between two closely related all-atom statistical energy
220 functions. *Protein Science* 17, 1212–1219.

221