1  **<u>Foldy: a web application for interactive protein structure analysis</u>**

2

3  **<u>Authors:</u>**

4  Jacob B. Roberts*[1,2,3], Alberto A. Nava*[1,2,4], Allison N. Pearson[1,2,5], Matthew R. Incha[1,2,5], Luis E.

5  Valencia[1,2,4], Melody Ma[6], Abhay Rao[3], Jay D. Keasling[1,2,3,4,7,8]

6

7  * Co-first authors

8   1. Joint BioEnergy Institute, Lawrence Berkeley National Laboratory, Emeryville, CA 94608, USA

9   2. Biological Systems and Engineering, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA

10  3. Department of Bioengineering, University of California, Berkeley, Berkeley, CA 94720, USA

11  4. Department of Chemical and Biomolecular Engineering, University of California, Berkeley, Berkeley, CA 94720, USA

12  5. Department of Plant and Microbial Biology, University of California, Berkeley, Berkeley, CA 94720, USA

13  6. Department of Molecular and Cell Biology, University of California, Berkeley, Berkeley, CA 94720, USA

14  7. Center for Synthetic Biochemistry, Shenzhen Institutes for Advanced Technologies, Shenzhen 518055, P.R. China

15  8. The Novo Nordisk Foundation Center for Biosustainability, Technical University Denmark, Kemitorvet, Building 220,

16     Kongens Lyngby 2800, Denmark

**Abstract**

17

18  Foldy is a cloud-based application that allows non-expert scientists to easily access and utilize

19  advanced AI-based structural biology tools, including AlphaFold and DiffDock. Built on

20  Kubernetes, it can be deployed by universities, departments, and labs without requiring hardware

21  resources, but can also be configured to utilize available computers. Foldy enables scientists to

22  predict the structure of proteins and complexes up to 3000 amino acids, visualize Pfam

23  annotations, and dock ligands with AutoDock Vina and DiffDock.

24

25  Our manuscript describes the user interface and deployment considerations of Foldy, as well as

26  some of our applications. By democratizing access to sophisticated AI-based tools, Foldy can

27  facilitate life science research and promote the wider adoption of structural bioinformatics tools.

28  Our work demonstrates that even the most advanced tools can be made accessible to a broad

29  audience through user-friendly platforms like Foldy, and we believe it will be a valuable resource

30  for researchers across scientific disciplines. The public structures available on the Lawrence

31  Berkeley Labs Foldy deployment can be viewed at https://foldy.lbl.gov.

32

33  **Author Summary**

34  Foldy is a cloud-based application that enables scientists to use AI-based structural biology tools

35  such as AlphaFold and DiffDock without software expertise. Built on Kubernetes, it can be set up

36  by universities, departments, and labs with no need for hardware resources. Foldy can predict

37  the structure of proteins and complexes up to 3000 amino acids, visualize Pfam annotations, and

38  dock ligands with AutoDock Vina and DiffDock. Our public structures can be viewed at

39  https://foldy.lbl.gov.

40

41  Our manuscript highlights the user interface, deployment considerations, and product applications

42  of Foldy. It's an accessible solution for researchers who are not software experts and can handle

43    the traffic of thousands of users and hundreds of thousands of protein structures and docked

44    ligands. This makes advanced AI-based tools more widely available, paving the way for

45    accelerating life science research.

46

47    By developing an easy-to-use platform, our work demonstrates that even the most sophisticated

48    AI-based tools can be made accessible to a wide audience. Foldy enables more scientists to draw

49    from the rapidly growing field of structural biology, making it a valuable tool for researchers across

50    scientific disciplines. We look forward to its adoption by the scientific community.

51

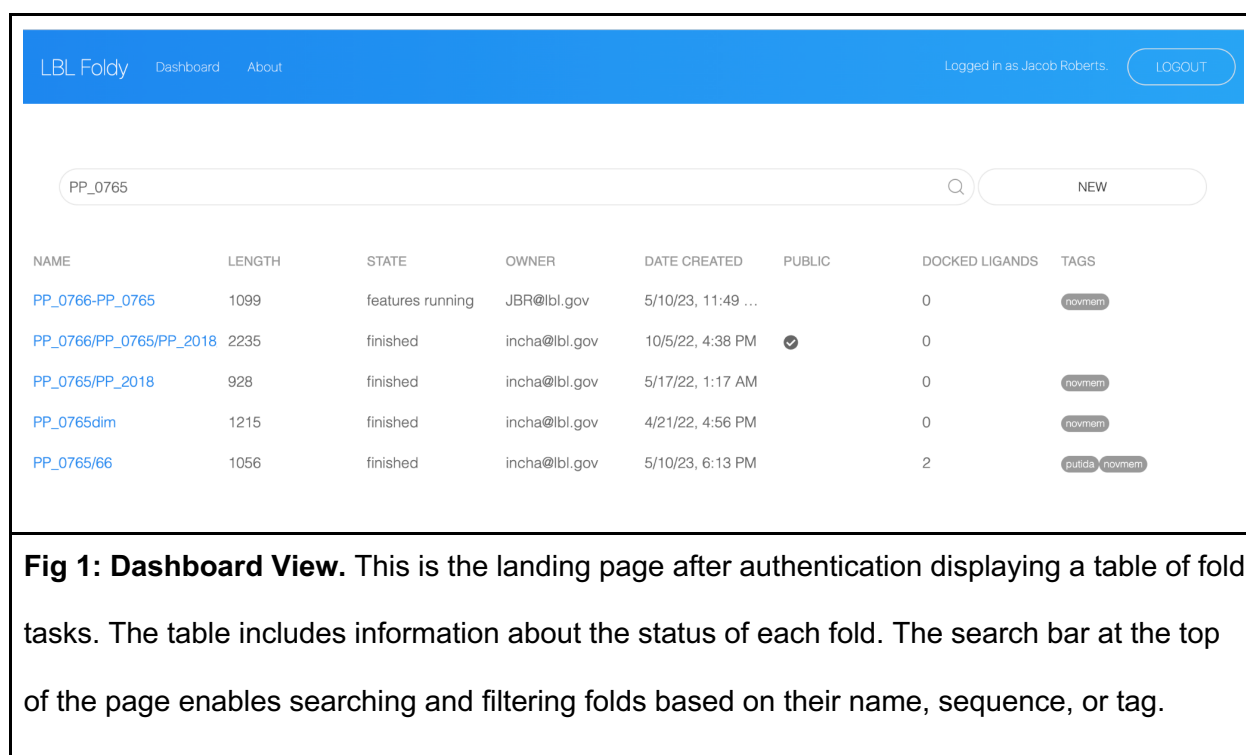52    **Introduction**

53    Recent advances in machine learning have led to the development of highly accurate protein

54    structure prediction methods [1–4], but their adoption has largely been limited to computational

55    biologists. These methods have produced impressive results in numerous applications including

56    *de novo* protein design [5] and protein-protein interaction screening [6]. However, the steep

57    requirements for storage space, GPU processing power, and RAM make the use of these tools

58    difficult for many end users. Nvidia created the BioNeMo service to make AI more accessible to

59    life science researchers, but it is a private implementation and currently only available to a few

60    biotechnology companies. Programs such as ColabFold [7] and AlphaFold-Colab [8] have been

61    developed to meet this need by providing custom Google Colaboratory Jupyter notebooks which

62    utilize free compute resources hosted by Google Cloud. These Jupyter notebooks provide an

63    interactive mode of using AlphaFold without the need for any complex installation or configuration.

64    However, there are several limitations to these notebooks including session timeouts, limited GPU

65    power, and limited batch processing capabilities. These issues are exacerbated when large

66    proteins (>1000 amino acids) are modeled which have higher resource demands. Ultimately,

67    these limitations mean that scaling up to more than a handful of structures is prohibitively difficult.

68

69  Here we present Foldy, an easy-to-deploy and easy-to-use modern web app for folding a protein

70  (AlphaFold[1]), predicting domain annotations (Pfam[9]), and docking small molecule ligands

71  (AutoDock Vina[10] or DiffDock[11]). Its primary objectives are to provide an intuitive interface,

72  facilitate deployment for IT administrators, and enable prediction tasks  for tens to thousands of

73  users per instance. The integrated tools within Foldy facilitate a seamless transition between

74  protein structure prediction and downstream analysis. It can be rapidly deployed in a cloud

75  environment, and also offers the possibility of a hybrid deployment utilizing local computational

76  resources. The design of the Foldy architecture aims to lower barriers to entry for both end users

77  and institutions.

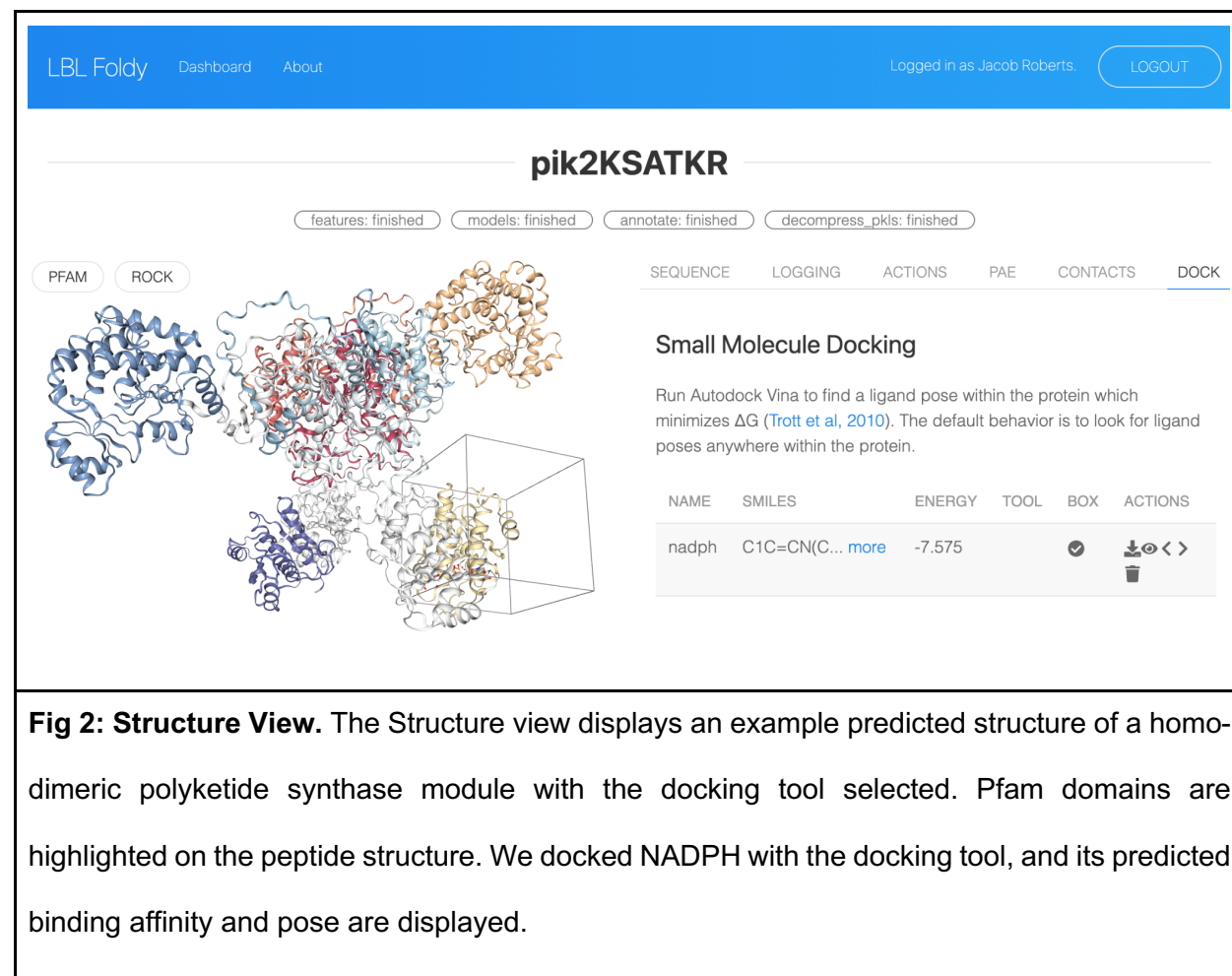78

79  **Results**

80  Interface

81  Foldy has four main views: the New Structure view, the Dashboard, the Tag view, and the

82  Structure view. The New Structure view (Fig S1) is where users can submit new structure

83  prediction tasks. At a minimum, users must provide an amino acid sequence and a name for the

84  structure. The Dashboard (Fig 1) serves as the app's landing page, providing access to all other

85  pages. By default, the Dashboard displays a table of the user's structures, but a search bar allows

86  users to filter structures by name, user, protein sequence, or tag. The Tag view (Fig S2) displays

87  all structures with a particular tag, and exposes bulk tasks such as downloading structures or

88  docking small molecule ligands.

89

| NAME | LENGTH | STATE | OWNER | DATE CREATED | PUBLIC | DOCKED LIGANDS | TAGS |
|---|---|---|---|---|---|---|---|
| PP_0766-PP_0765 | 1099 | features running | JBR@lbl.gov | 5/10/23, 11:49 … | | 0 | novmem |
| PP_0766/PP_0765/PP_2018 | 2235 | finished | incha@lbl.gov | 10/5/22, 4:38 PM | ✓ | 0 | |
| PP_0765/PP_2018 | 928 | finished | incha@lbl.gov | 5/17/22, 1:17 AM | | 0 | novmem |
| PP_0765dim | 1215 | finished | incha@lbl.gov | 4/21/22, 4:56 PM | | 0 | novmem |
| PP_0765/66 | 1056 | finished | incha@lbl.gov | 5/10/23, 6:13 PM | | 2 | putida novmem |

**Fig 1: Dashboard View.** This is the landing page after authentication displaying a table of fold tasks. The table includes information about the status of each fold. The search bar at the top of the page enables searching and filtering folds based on their name, sequence, or tag.

90

91   There are two types of users: editors and viewers. Editors have full read and write access. Viewers

92   are any user with a Google account, and are allowed view-only access to structures which have

93   been explicitly marked "public" and their associated data (logs, docking runs). Users are

94   authenticated by their Gmail account, and user types are flag controlled. You can view our public

95   structures at https://foldy.lbl.gov.

96

97   The Structure view has two columns: the predicted structure is on the left and a tool panel is on

98   the right (Fig 2). By default the tool panel displays the amino acid sequence (Fig 3A), and Pfam

99   domain annotations can be overlaid on both the structure and sequence(s) (Fig 2 left, Fig 3A). A

100  number of actions are available to users through the tabs in the tool panel. For example, users

101  can predict residue interactions and complex formation using contact probability maps (Fig 3E)

102  [6]. Users can segment proteins into domains and predict inter-domain flexibility using the

103  Predicted Alignment Error (PAE, Fig 3D) [3]. Additionally, users can dock small molecule ligands

5

104    with AutoDock Vina or DiffDock by specifying the SMILES string and optionally a bounding box

105    around a residue (Fig 3F) [10,11].

106



**Fig 2: Structure View.** The Structure view displays an example predicted structure of a homo-dimeric polyketide synthase module with the docking tool selected. Pfam domains are highlighted on the peptide structure. We docked NADPH with the docking tool, and its predicted binding affinity and pose are displayed.

107

108



7

**Fig 3: Structure View Toolbar Tabs.** *A. Sequence*: this tab displays the peptide sequences, optionally annotated with domain predictions, such as the pfam annotations displayed here. *B. Logging*: this tab displays real-time logs of each sub-task within a protein structure prediction task. *C. Actions*: this tab allows a user to download the top ranked pdb structure file of a given fold as well as all the intermediate files generated by the fold pipeline *D. PAE*: this tab shows the predicted alignment error heatmap, measured in Ångstroms, between each residue and between each chain. The PAE is derived from the AlphaFold model. *E. Contact Probability*: this tab shows the contact probability heatmap between each residue and between each chain. The contact probabilities are derived from the AlphaFold model. *F. Docking*: this tab shows a table with any docking results as well as a submission form for new docking tasks. The docking submission form requires a ligand name and SMILES string, and optionally a bounding box can be specified of a particular radius around a particular residue (Vina only).

109

110    Facile Deployment

111    A local version of Foldy can be set up in seconds with Docker Compose, but none of the analytical

112    tools are supported locally. The full Foldy can be deployed in the cloud in a matter of hours by an

113    IT administrator. It requires a few manually provisioned resources: a web domain for the app, a

114    static IP address, a database, and a Kubernetes cluster. The remainder of the app, including the

115    frontend, backend, and compute workers, are all managed by a single Helm chart. Helm is a

116    command line tool for managing Kubernetes configuration, and with one command can deploy all

117    the necessary Kubernetes resources. The infrastructure schematic is illustrated in Figure S3, and

118    the Helm chart with corresponding step-by-step deployment instructions are available at

119    https://github.com/JBEI/foldy.

120

121     <u>Scalability</u>

122     By default, Foldy uses cloud compute for all tasks, meaning any lab or institution, regardless of

123     their hardware, is able to set up Foldy. In this setup, jobs are run on ephemeral machines in the

124     cloud that are spawned when work tasks are queued and deleted when the work is complete. The

125     work tasks are tracked in a queue (implemented with RedisQueue), and the worker machines are

126     automatically created and destroyed by Kubernetes (implemented with Prometheus and KEDA).

127     Importantly, each worker machine can be provisioned with high-memory GPUs, enabling the

128     prediction of large protein structures.

129

130     Institutions with access to their own compute resources, including compute clusters, can run the

131     worker threads on their own machines. To run a worker on a local compute resource which

132     supports docker, one can run the "worker" docker image on the machines with the appropriate

133     flags, and set up tunnels to the cloud databases. To use local compute resources which don't

134     support docker, one can create bash scripts which execute each tool. For example, to run

135     AlphaFold jobs on a university cluster which does not support docker, one must create a variant

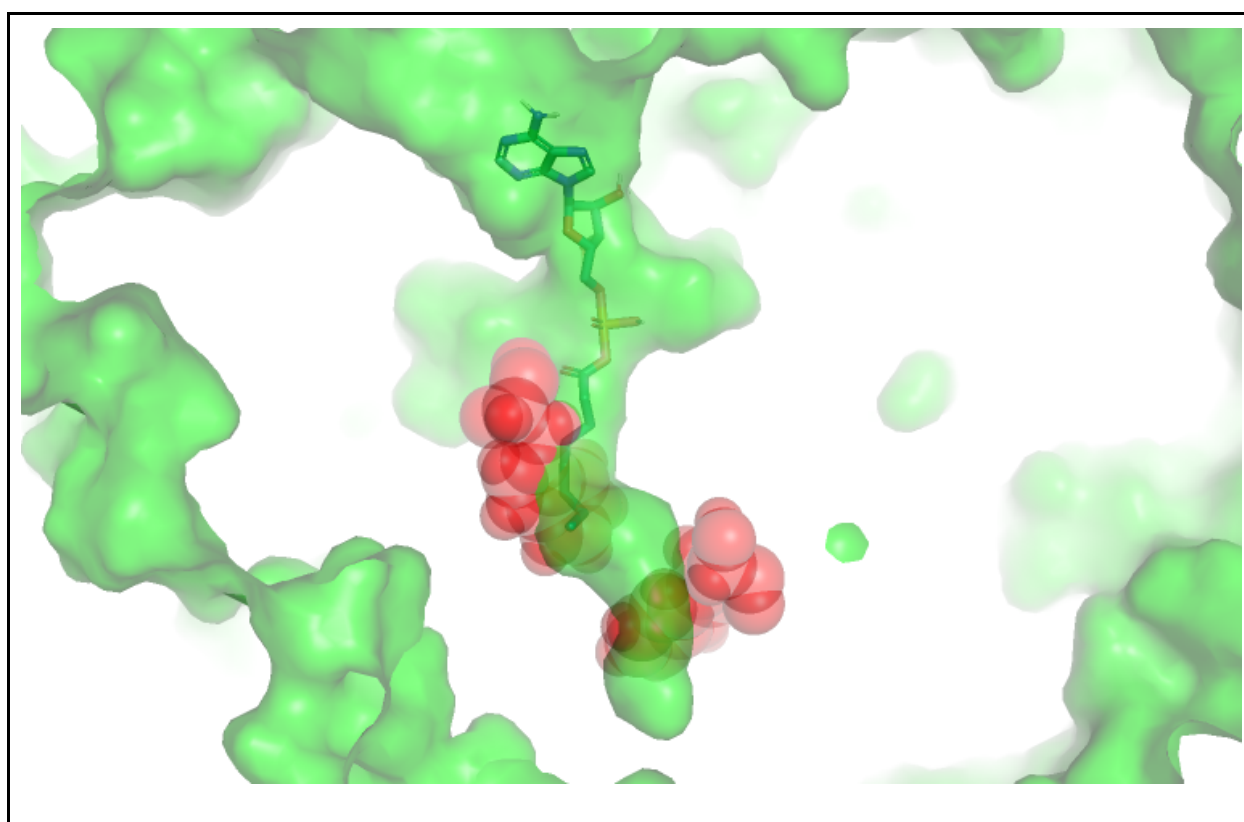136     of "worker/run_alphafold.sh" which invokes the local AlphaFold installation.

137

138     **<u>Discussion</u>**

139     This tool greatly facilitates research because it makes complex tools accessible. The LBL Foldy

140     instance (https://foldy.lbl.gov) has been used by 55 researchers across 6 labs, to predict 4802

141     structures and dock 2754 ligands. It has been used in over a dozen projects, three of which are

142     case-studied below.

143

144     <u>Foldy Case Studies</u>

145     One researcher was interested in changing the substrate preference of a long-chain fatty acyl-

146     AMP ligase (FAAL) from long-chain to medium-chain. They used Foldy to predict the wild type

147     protein's structure, and used AutoDock Vina to dock octanoyl-AMP in the active site tunnel. They

148     considered two point mutations to tryptophan which they hypothesized would obstruct the

149     substrate tunnel and shorten chain-length preference of the FAAL. They found from the structure

150     that one tryptophan point mutation would occlude even C8 compounds, while the other would
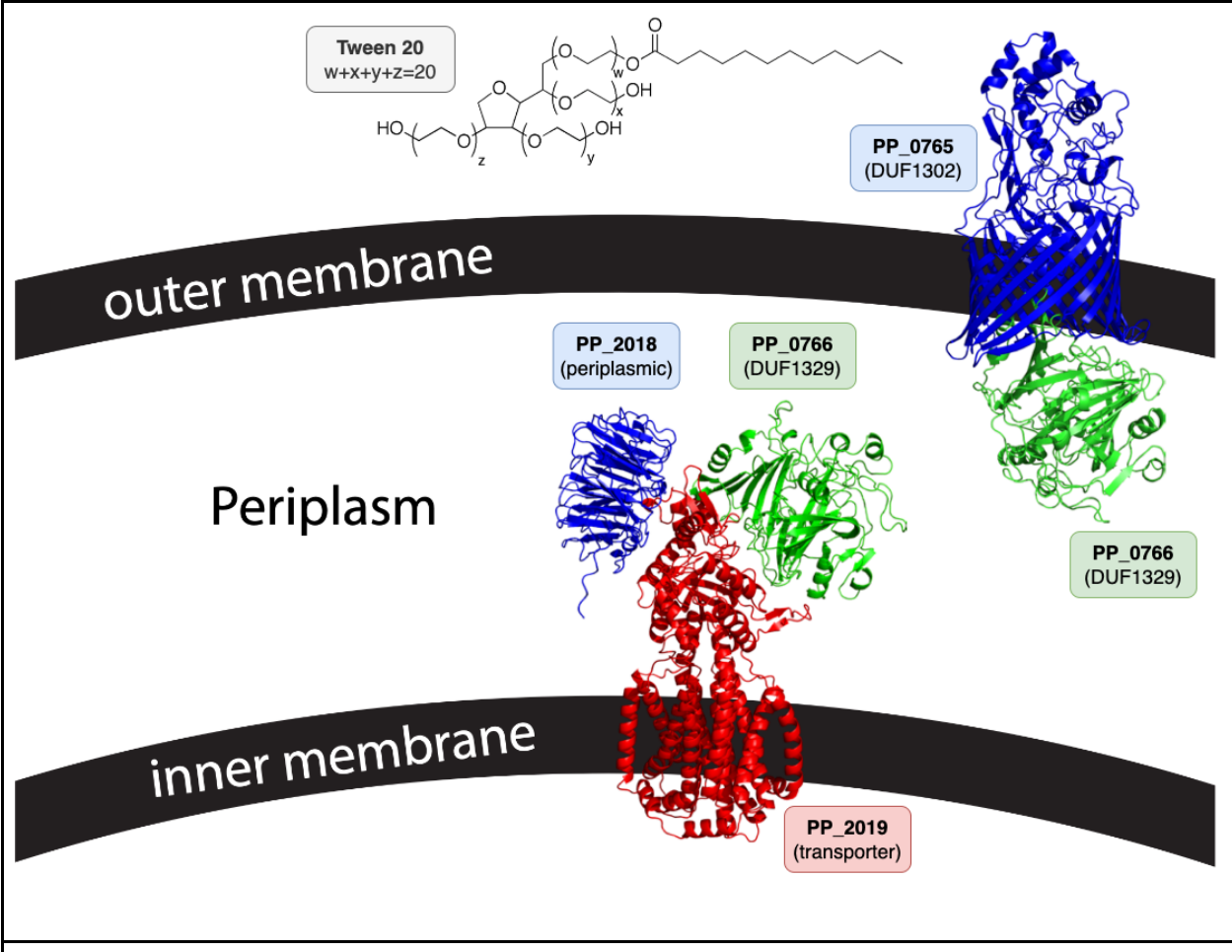
151     occlude C10 or larger (Fig 4).

152



**Fig 4: A fatty acyl-AMP ligase with substrate docked.** The two red sections represent two point mutations under consideration: the location of the docked ligand shows that the upper point mutation would occlude the pictured C8 ligand, while the lower point mutation would occlude ligands C10 or larger.

153

154     Foldy was used to evaluate AlphaFold Multimer's ability to predict chimeric polyketide synthase

155     production. Nava et al (in review) evaluated whether a chimeric PKS's predicted structure is

156    indicative of its production titer. A total of 144 interacting KS-AT / ACP pairs from a seminal paper

157    about module swaps [12] were co-folded with AlphaFold Multimer in Foldy, and the protein contact

158    probability map was used as a proxy for likelihood of protein interaction, as described by

159    Humphrey's et al[6]. The authors found the predicted structures informative and worth more

160    investigation.

161

162    Metabolic engineers used Foldy to predict the function of dozens of domains of unknown function

163    (DUFs) in *P. putida*, including DUF1302 and DUF1329. These two DUFs which were previously

164    suspected of being involved in hydrolase activity [13] may actually be involved in substrate

165    transport. DUFs with pfam IDs DUF1302 and DUF1329, represented in *P. putida* by PP_0765

166    and PP_0766, have no predicted function, but prior RB-TnSeq experiments show their function

167    correlates with both a periplasmic protein (PP_2018) and a multidrug efflux transporter (PP_2019)

168    [13]. Combinatorial protein-protein docking sims, done with AlphaFold in Foldy, predicted two

169    complexes. First, the two DUFs seem to interact with one another: PP_0765, which appears to

170    be a membrane bound beta barrel, is predicted to form a complex with PP_0766 (Fig 5, top).

171    Second, PP_0766 is predicted to form a heterotrimer with both PP_2018 and PP_2019 (Fig 5,

172    bottom). This may indicate that these proteins are not in fact involved in hydrolase activity as

173    previously reported[13] but rather are components of a novel transport system. Future work

174    should interrogate this system in greater detail to determine the directionality of transport and *in*

175    *vivo* function

176

11

**Fig 5: Putative complex formation of two domains of unknown function in *P. putida*.**

Two DUFs in *P. putida*, previously hypothesized to have hydrolase activity may actually be involved in substrate transport: PP_0765 (DUF1302, top blue) and PP_0766 (DUF1329, top green, bottom green). PP_0765 has the characteristic beta-barrel of a membrane protein, and shows high likelihood of forming a complex with PP_0766 (top). Additionally, PP_0766 is predicted to form a heterotrimer with PP_2018 (bottom blue) and PP_2019 (bottom red).

177

178

179    Conclusion

180    Foldy is easier to use than other AlphaFold implementations, and addresses some of their

181    limitations [1,7]. User experience studies indicate that small improvements in the enjoyability of a

182    tool may have significant effects on tool use [14]. Foldy will increase biologists' productivity by

183    increasing the adoption of computational structure tools. The adoption of Foldy by institutions,

184    including those without large compute clusters or GPUs, will make high-accuracy protein structure

185    prediction more accessible.

186

187    **Availability and Future Developments**

188    Website built on Kubernetes, with the Helm chart and Dockerfiles freely available at

189    https://github.com/JBEI/foldy. One can access the "public" structures in our group's deployment

190    at foldy.lbl.gov.

191

192    **Acknowledgements**

208

209 J.D.K. has financial interests in Amyris, Ansa Biotechnologies, Apertor Pharma, Berkeley Yeast,

210 Demetrix, Lygos, Napigen, ResVita Bio, and Zero Acre Farms.

211

212 **<u>References</u>**

213 1. Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O, et al. Highly accurate

214     protein structure prediction with AlphaFold. Nature. 2021;1–11.

215 2. Baek M, DiMaio F, Anishchenko I, Dauparas J, Ovchinnikov S, Lee GR, et al. Accurate

216     prediction of protein structures and interactions using a three-track neural network.

217     Science. 2021;eabj8754.

218 3. Evans R, O'Neill M, Pritzel A, Antropova N, Senior A, Green T, et al. Protein complex

219     prediction with AlphaFold-Multimer. bioRxiv. 2021;2021.10.04.463034.

220 4. Lin Z, Akin H, Rao R, Hie B, Zhu Z, Lu W, et al. Evolutionary-scale prediction of atomic

221     level protein structure with a language model [Internet]. bioRxiv. 2022 [cited 2022 Nov 2]. p.

222     2022.07.20.500902. Available from:

223     https://www.biorxiv.org/content/10.1101/2022.07.20.500902v2

224 5. Wang J, Lisanza S, Juergens D, Tischer D, Watson JL, Castro KM, et al. Scaffolding

225     protein functional sites using deep learning. Science. 2022 Jul 22;377(6604):387–94.

226 6. Humphreys IR, Pei J, Baek M, Krishnakumar A, Anishchenko I, Ovchinnikov S, et al.

227     Computed structures of core eukaryotic protein complexes. Science.

228     2021;374(6573):eabm4805.

229 7. Mirdita M, Schütze K, Moriwaki Y, Heo L, Ovchinnikov S, Steinegger M. ColabFold: making

230     protein folding accessible to all. Nat Methods. 2022 Jun;19(6):679–82.

231   8.   Tunyasuvunakool K, Adler J, Wu Z, Green T, Zielinski M, Žídek A, et al. Highly accurate

232        protein structure prediction for the human proteome. Nature. 2021;596(7873):590–6.

233   9.   Mistry J, Chuguransky S, Williams L, Qureshi M, Salazar GA, Sonnhammer ELL, et al.

234        Pfam: The protein families database in 2021. Nucleic Acids Res. 2021 Jan 8;49(D1):D412–

235        9.

236   10.  Eberhardt J, Santos-Martins D, Tillack AF, Forli S. AutoDock Vina 1.2.0: New Docking

237        Methods, Expanded Force Field, and Python Bindings. J Chem Inf Model. 2021 Aug

238        23;61(8):3891–8.

239   11.  Corso G, Stärk H, Jing B, Barzilay R, Jaakkola T. DiffDock: Diffusion Steps, Twists, and

240        Turns for Molecular Docking [Internet]. arXiv [q-bio.BM]. 2022. Available from:

241        http://arxiv.org/abs/2210.01776

242   12.  Menzella HG, Carney JR, Santi DV. Rational Design and Assembly of Synthetic Trimodular

243        Polyketide Synthases. Chem Biol. 02/2007;14(2):143–51.

244   13.  Thompson MG, Incha MR, Pearson AN, Schmidt M, Sharpless WA, Eiben CB, et al. Fatty

245        Acid and Alcohol Metabolism in Pseudomonas putida: Functional Analysis Using Random

246        Barcode Transposon Sequencing. Appl Environ Microbiol [Internet]. 2020 Oct 15;86(21).

247        Available from: http://dx.doi.org/10.1128/AEM.01665-20

248   14.  Diefenbach S, Hassenzahl M. The dilemma of the hedonic – Appreciated, but hard to

249        justify. Interact Comput. 2011 Sep;23(5):461–72.

250