

ColabFold - Making protein folding accessible to all

Milot Mirdita,^{1,*} Konstantin Schütze,² Yoshitaka Moriwaki,^{3,4} Lim Heo,⁵ Sergey Ovchinnikov,^{6,7,*} and Martin Steinegger^{2,8,*}

ColabFold offers accelerated protein structure and complex predictions by combining the fast homology search of MMseqs2 with AlphaFold2 or RoseTTAFold. ColabFold's 40–60× faster search and optimized model use allows predicting close to a thousand structures per day on a server with one GPU. Coupled with Google Colaboratory, ColabFold becomes a free and accessible platform for protein folding. ColabFold is open-source software available at github.com/sokrypton/ColabFold. Its novel environmental databases are available at colabfold.mmseqs.com

Contact: milot.mirdita@mpibpc.mpg.de, so@fas.harvard.edu, martin.steinegger@snu.ac.kr

1 Predicting the three-dimensional structure of a protein from
2 its sequence alone remains an unsolved problem. However,
3 by exploiting the information in multiple sequence alignments
4 (MSAs) of related proteins as raw input features for end-to-
5 end training, AlphaFold2 [1] was able to predict the 3D atomic
6 coordinates of folded protein structures at a median GDT-TS
7 of 92.4% in the latest CASP14 [2] competition. The accuracy
8 of many of the predicted structures was within the error mar-
9 gin of experimental structure determination methods. Many
10 ideas of AlphaFold2 were independently reproduced and im-
11 plemented in RoseTTAFold [3]. Additionally to single chain
12 predictions, RoseTTAFold was shown to model protein com-
13 plexes. Evans *et al.* [4] released AlphaFold-multimer, a re-
14 fined version of AlphaFold2 for complex prediction. Thus,
15 two highly accurate open-source prediction methods are now
16 publicly available.

17 In order to leverage the power of these methods re-
18 searchers require powerful compute-capabilities. First, to
19 build diverse MSAs, large collections of protein sequences
20 from public reference [5] and environmental [1, 6] databases
21 are searched using the most sensitive homology detection
22 methods HMMer [7] and HHblits [8]. These environmental
23 databases contain billions of proteins extracted from metage-
24 nomic and -transcriptomic experiments, which often comple-
25 ment databases dominated by isolate genomes. Due to their
26 large size searches can take up to hours for a single protein,
27 while requiring over two terabyte of storage space alone. Sec-
28 ond, to execute the deep neural networks GPUs with a large
29 amount of GPU RAM are required even for relatively common
30 protein sizes of ~1000 residues. Though, for these the MSA
31 generation dominates the overall run-time.

32 To enable researchers without these resources to use Al-
33 phaFold2, independent solutions based on Google Colabora-
34 tory were developed. Colaboratory is a proprietary version
35 of Jupyter Notebook hosted by Google. It is accessible for
36 free to logged-in users and includes access to powerful GPUs.
37 Tunyasuvunakool *et al.* [9] developed an AlphaFold2 Jupyter

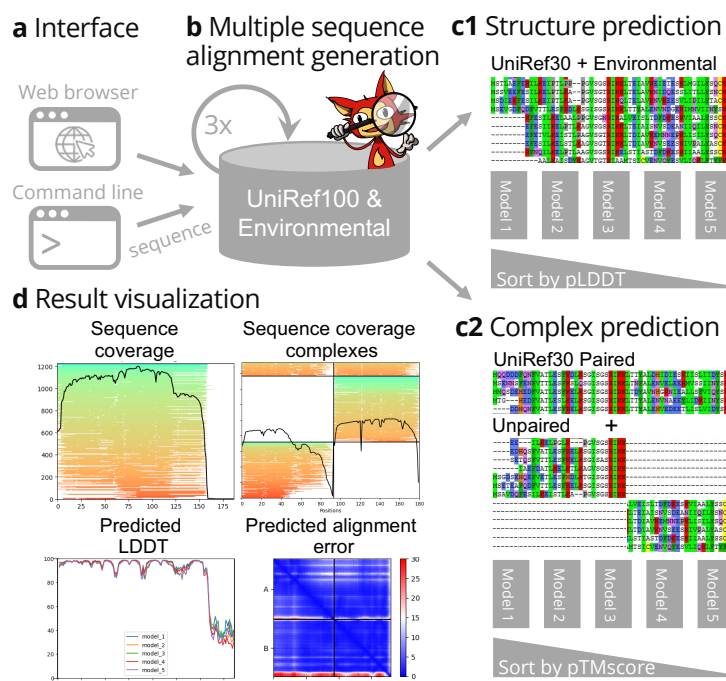


FIG. 1. (a) ColabFold has a web and a command line interface, that (b) send FASTA input sequence(s) to a MMseqs2 server searching two databases UniRef100 and a database of environmental sequences with three profile-search iterations each. The second database is searched using a sequence-profile generated from the UniRef100 search as input. The server generates two MSAs in A3M format containing all detected sequences. (c1) For single structure predictions we filter both A3Ms using a diversity aware filter and return this to be provided as the MSA input feature to the AlphaFold2 models. (c2) For complex prediction we pair the top hits within the same species to resolve the inter-complex contacts and additionally add two unpaired MSAs (same to c1) to guide the structure prediction. (d) To help researchers judge the prediction quality we visualize MSA depth and diversity and show the AlphaFold2 confidence measures (pLDDT and PAE).

38 Notebook for Google Colaboratory (referred to as AlphaFold-
39 Colab), where the input MSA is built by searching with HM-
40 Mer against a clustered UniProt and an eight-fold reduced en-
41 vironmental databases. Resulting in less accurate predictions,
42 while still requiring long search times.

43 Here, we present ColabFold, a fast and easy to use software
44 for protein structure and homo- and heteromer complex pre-
45 diction, for use as a Jupyter Notebook inside Google Colabora-
46 tory, on researchers' local computers as a notebook or through
47 a command line interface. ColabFold speed-ups the predic-
48 tion by replacing AlphaFold2's homology search with a 40-60

¹Quantitative and Computational Biology, Max Planck Institute for Biophysical Chemistry, Göttingen, Germany. ²School of Biological Sciences, Seoul National University, Seoul, South Korea. ³Department of Biotechnology, Graduate School of Agricultural and Life Sciences, The University of Tokyo, Tokyo, Japan. ⁴Collaborative Research Institute for Innovative Microbiology, The University of Tokyo, Tokyo, Japan. ⁵Department of Biochemistry and Molecular Biology, Michigan State University, East Lansing, MI 48824, USA. ⁶JHDSF Program, Harvard University, Cambridge, MA 02138, USA. ⁷FAS Division of Science, Harvard University, Cambridge, MA 02138, USA. ⁸Artificial Intelligence Institute, Seoul National University, Seoul, South Korea * These authors contributed equally and are ordered alphabetically.

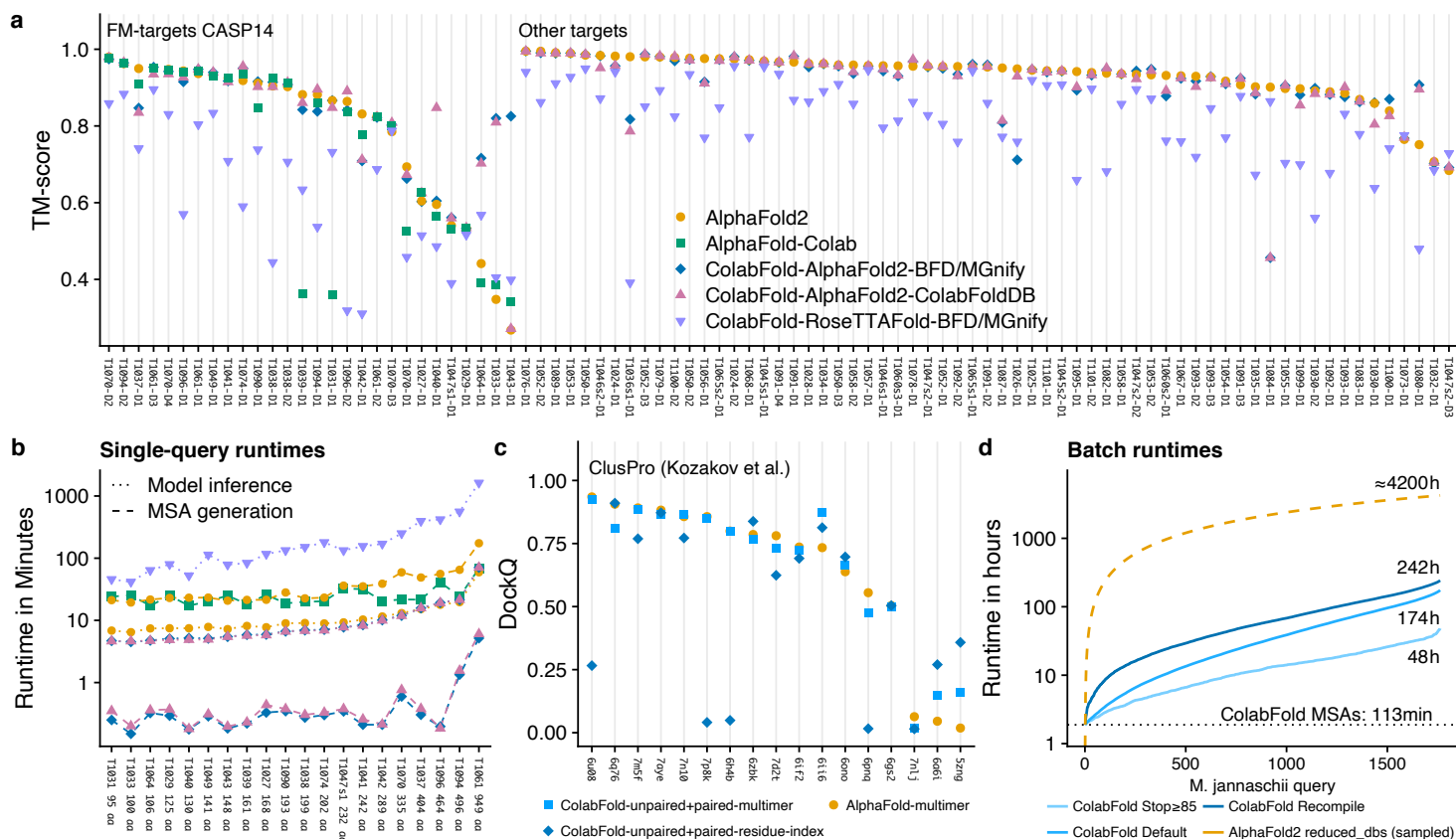


FIG. 2. (a) Structure prediction comparison of AlphaFold2 (yellow), AlphaFold-Colab (green) and ColabFold-AlphaFold2 with BFD/MGnify (blue) and with the ColabFoldDB (magenta), and ColabFold-RoseTTAFold with BFD/MGnify (purple) using predictions of 91 domains of 65 CASP14 targets. The 28 domains from the 20 free-modeling (FM) targets are shown first. FM targets were used to optimize MMseqs2 search parameters. Each target was evaluated for each individual domain (in total 91 domains). (b) MSA generation and model inference times for each CASP14 FM target sorted by protein length (same colors as before). Blue shows MSA runtimes for ColabFold-AlphaFold2-BFD/MGnify and ColabFold-RoseTTAFold-BFD/MGnify. (c) Comparison of ColabFold complex predictions in residue-index- (dark blue) and AlphaFold-multimer (light blue) mode, and to AlphaFold-multimer (yellow). (d) Runtime of colabfold_batch proteome prediction at three optimization levels: (dark blue) Always recompile, (blue) default, (light blue) stop model/recycle evaluation after first prediction with a pLDDT of ≥ 85 . Extrapolated line based on 50 AlphaFold2 predictions shown in yellow.

49 times faster MMseqs2 [10, 11] search. It additionally imple-
 50 ments speed-ups for batch predictions of structures by avoid-
 51 ing recompilation and adding early stop criteria. ColabFold’s
 52 batch mode with early stopping can compute the proteome of
 53 *Methanocaldococcus jannaschii* in 48 h on a consumer GPU –
 54 a ~ 90 times speedup over AlphaFold2. We show that Colab-
 55 Fold outperforms AlphaFold-Colab and matches AlphaFold2
 56 on CASP14 targets and also matches AlphaFold-multimer on
 57 the ClusPro [4, 12] dataset in prediction quality.

58 ColabFold (Fig. 1) consists of three parts: (1) An MMseqs2
 59 based homology search server to build diverse MSAs and to
 60 find templates. The server efficiently aligns input sequence(s)
 61 against the UniRef100, the PDB70 and an environmental se-
 62 quence set. (2) A Python library that communicates with the
 63 MMseqs2 search server, prepares the input features for (single
 64 or complex) structure inference, and visualizes of results. This
 65 library also implements a command line interface. (3) Jupyter
 66 notebooks for basic, advanced and batch use (Methods “Col-
 67 labFold notebooks”) using the Python library.

68 In ColabFold we replace the sensitive search methods HM-
 69 Mer and HHblits by MMseqs2. We optimized the MSA genera-
 70 tion by MMseqs2 to have the following three properties: (1)
 71 MSA generation should be fast. (2) The MSA has to capture
 72 diversity well and (3) it has to be small enough to run on
 73 computers with limited RAM. Reducing the memory require-
 74 ment is especially helpful in Google Colaboratory where the
 75 provided system is selected from a pool with widely differing
 76 capabilities. While (1) is achieved through the fast MMseqs2
 77 prefilter for (2 and 3) we developed a search workflow to maxi-
 78 mize sensitivity (Methods “MSA generation”) and a new filter
 79 that samples the sequence space evenly (Methods “New diver-
 80 sity aware filter” and **Supplementary Fig. 1**). Prediction
 81 quality highly depends on the input MSA. However, often an
 82 MSA with only a few (~ 30) sufficiently diverse sequences is
 83 enough to produce high quality predictions (see Jumper et al.,
 84 **Fig. 5a**).

85 Additionally, we combined the BFD and MGnify databases
 86 that are used in AlphaFold2 by HHblits and HMMer respec-

87 tively into a combined redundancy reduced version we refer to
88 as BFD/MGNify (Methods “Reducing size of BFD/MGNify”).
89 The environmental search database presented an opportunity
90 to improve structure predictions of non-bacterial sequences,
91 as e.g., eukaryotic protein diversity is not well represented in
92 the BFD and MGNify databases. Limitations in assembly and
93 gene calling due to complex intron/exon structures result in
94 under representation in reference databases. We therefore ex-
95 tended the BFD/MGNify with additional metagenomic protein
96 catalogues containing eukaryotic proteins [13, 14, 15], phage
97 catalogues [16, 17] and an updated version of MetaClust [18].
98 We refer to this database as ColabFoldDB (Methods “Colab-
99 FoldDB”). In **Supplementary Fig. 2** we show that the Colab-
100 FoldDB in comparison to the BFD/MGNify produces more
101 diverse MSAs for PFAM [19] domains with < 30 members.

102 To compare the accuracy of predicted structures we
103 compared AlphaFold2 (default settings with templates),
104 AlphaFold-Colab (no templates), ColabFold-RoseTTAFold-
105 BFD/MGNify, ColabFold-AlphaFold2-BFD/MGNify and
106 ColabFold-AlphaFold2-ColabFoldDB on TM-scores for all
107 targets from the CASP14 competition (**Fig. 2a**). All three
108 ColabFold modes were executed without templates. We show
109 the targets split by free modeling (FM) on the left and the
110 remaining ones on the right, since we used the FM-targets for
111 optimization of search workflow parameters.

112 The mean TM-scores for the FM targets are 0.826,
113 0.818, 0.79, 0.744 and 0.62 for ColabFold-AlphaFold2-
114 BFD/MGNify, ColabFold-AlphaFold2-ColabFoldDB, Al-
115 phaFold2, AlphaFold-Colab and ColabFold-RoseTTAFold-
116 BFD/MGNify respectively. Over all CASP14 targets the
117 TM-scores are 0.887, 0.886, 0.888 and 0.754 for the respective
118 methods, excluding AlphaFold-Colab as it cannot be used
119 stand-alone.

120 ColabFold could not predict T1084 well as MMseqs2 sup-
121 presses all databases hits as false positives due to its amino
122 acid composition filter and masking procedure. If these filters
123 are deactivated T1084 can be predicted with an TM-score of
124 0.872 (**Supplementary Fig. 3**). **Supplementary Table 1**
125 contains a list of further targets where ColabFold differed sig-
126 nificantly from AlphaFold2.

127 ColabFold is on average 5x faster for single predictions than
128 AlphaFold2 and AlphaFold-Colab, when taking both MSA
129 generation (**Fig. 2b**) and model inference into account.

130 AlphaFold2 was initially released without capabilities to
131 model complexes. However, we found that by combining two
132 sequences with a glycine linker [20] it could often successfully
133 model complexes. Shortly afterwards, Baek [21] found that in-
134 crementing the model-internal residue index - the method that
135 was used in RoseTTAFold - could also be used in AlphaFold2.

136 For high quality predictions it was shown that sequences
137 should be provided in paired-form to AlphaFold2 [22]. We im-
138 plemented a similar pairing procedure (Methods “MSA pair-
139 ing for complex prediction”) and show the complex prediction
140 capabilities of ColabFold in **Fig. 2c**. ColabFold achieves the
141 highest accuracy in complex prediction on the ClusPro [4, 12]
142 dataset with the AlphaFold-multimer model, however, some
143 targets performed better using the residue-index mode.

144 **Fig. 3** shows two examples of ColabFold’s complex predic-

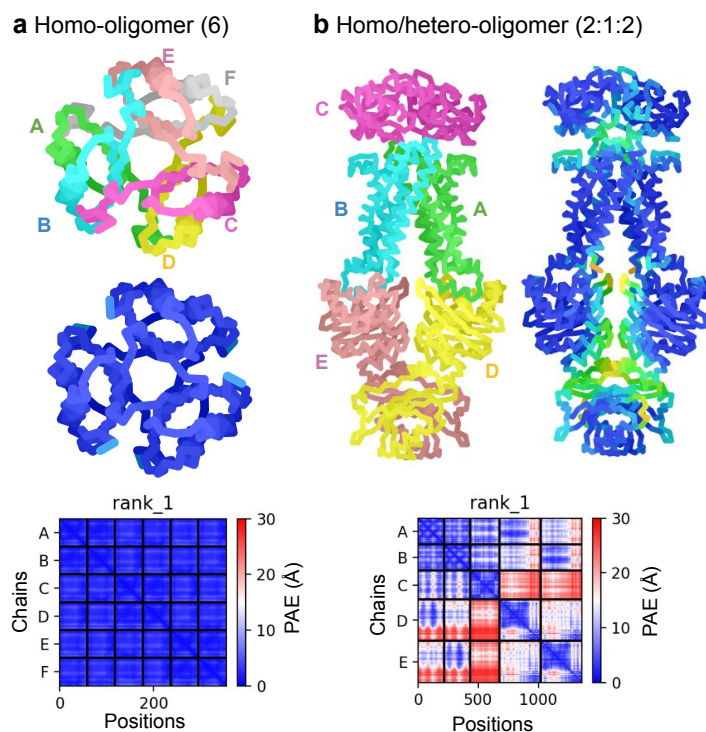


FIG. 3. Anecdotal examples showcasing the capabilities of advanced ColabFold features. (a) Setting the homo-oligomer setting to 6, allows modeling of the homo-6-mer structure of 4-Oxalocrotonate Tautomerase. Colored by chain (top), pLDDT (predicted Local Distance Difference Test, bottom). The inter PAE (Predicted Aligned Error) between chains is very low indicating a confident prediction. (b) Providing three different proteins with 2:1:2 homo-oligomer setting allows modeling a hetero-complex with mismatching symmetries of the D-methionine transport system.

145 tion capabilities: (a) shows a homo-six-mer and (b) shows
146 a D-methionine transport system composed of three different
147 proteins. For single structure prediction AlphaFold2 provides
148 a pLDDT measure to indicate the prediction quality. A high
149 pLDDT does not necessarily indicate a correct complex pre-
150 diction, though the inter-complex predicted alignment error
151 (PAE) helps to rank complexes. We visualize plots of PAE
152 and complex conformation to help users judge the prediction
153 quality of a complex. An example for heteromer complex pre-
154 diction is shown in **Supplementary Fig. 4** with its PAE plot.
155 Furthermore, ColabFold complexes were successfully used to
156 aid the cryo-EM structure determination of the 120 MDa hu-
157 man nucleopore complex [23].

158 In ColabFold we expose many internal parameters of Al-
159 phaFold2 to aid users to model difficult targets, such as the
160 recycle count (default 3). It controls the number of times
161 the prediction is repeatedly fed through the model. For dif-
162 ficult targets as well as for designed proteins without known
163 homologs additional recycling iterations can result in a high
164 quality prediction (**Supplementary Fig. 5**). Rerunning the
165 CASP14 benchmark using 12 recycles resulted in an improve-
166 ment of average TM-score from 0.887 to 0.898 (**Supplemen-
167 tary Fig. 6**). The largest improvement was in targets with
168 little MSA information.

169 To meet the demand for high throughput structure predic-

170 tion we introduced several features in ColabFold. (1) MSA
171 generation can be executed in batch-mode independently from
172 model batch-inference. (2) We compile only one of the five Al-
173 phaFold2 models and reuse weights. (3) We provide a batch
174 execution mode, that avoids recompilation for sequences of
175 similar length. (4) We implement early stop criteria, to avoid
176 running additional recycles or models if a sufficiently accurate
177 structure was already found. (5) We developed the command
178 line tool `colabfold_batch` to predict structures on local ma-
179 chines. All together, we show that the proteome of 1762 pro-
180 teins shorter than 1000 aa of *M. jannaschii* can be predicted in
181 48 h with early stopping at pLDDT of ≥ 85 on one Nvidia Titan
182 RTX (**Fig. 2d**), while sacrificing little-or-no prediction accu-
183 racy (Methods “Proteome Benchmark”). The average pLD-
184 DTs of AlphaFold2 and ColabFold Stop ≥ 85 were 89.75 and
185 88.78 in a subsampled set of 50 proteins.

186 ColabFold builds beyond the initial offerings of Alphafold2
187 by improving its sequence search, providing tools for model-
188 ing homo- and heteromer complexes, exposing advanced func-
189 tionality, expanding the environmental databases and enabling
190 large-scale batch prediction of protein structures – at a ~ 90
191 times speedup over AlphaFold2.

REFERENCES

- [1] Jumper, J. *et al. Nature* **596**, 583–589 (2021).
- [2] Kryshchuk, A. *et al. Proteins* **89**, 1607–1617 (2021).
- [3] Baek, M. *et al. Science* **373**, 871–876 (2021).
- [4] Evans, R. *et al. bioRxiv* 2021.10.04.463034 (2021).
- [5] UniProt Consortium. *Nucleic Acids Res.* **47**, D506–D515 (2019).
- [6] Mitchell, A. L. *et al. Nucleic Acids Res.* **48**, D570–D578 (2020).
- [7] Eddy, S. R. *PLoS Comput. Biol.* **7**, e1002195 (2011).
- [8] Steinegger, M. *et al. BMC Bioinform.* **20**, 473 (2019).
- [9] Tunyasuvunakool, K. *et al. Nature* **596**, 590–596 (2021).
- [10] Steinegger, M. & Söding, J. *Nat. Biotechnol.* **35**, 1026–1028 (2017).
- [11] Mirdita, M. *et al. Bioinformatics* **35**, 2856–2858 (2019).
- [12] Kozakov, D. *et al. Nat. Protoc.* **12**, 255–278 (2017).
- [13] Levy Karin, E. *et al. Microbiome* **8**, 48 (2020).
- [14] Delmont, T. O. *et al. bioRxiv* 2020.10.15.341214 (2020).
- [15] Alexander, H. *et al. bioRxiv* 2021.07.25.453713 (2021).
- [16] Nayfach, S. *et al. Nat. Microbiol.* **6**, 960–970 (2021).
- [17] Camarillo-Guerrero, L. F. *et al. Cell* **184**, 1098–1109.e9 (2021).
- [18] Steinegger, M. & Söding, J. *Nat. Commun.* **9**, 2542 (2018).
- [19] Mistry, J. *et al. Nucleic Acids Res.* **49** (2021).
- [20] Moriwaki, Y. AlphaFold2 can also predict heterocom-
plexes. all you have to do is input the two sequences
you want to predict and connect them with a long linker.
https://twitter.com/Ag_smith/status/1417063635000598528 (2021).
- [21] Baek, M. Adding a big enough number for “residue_index” feature is
enough to model hetero-complex using AlphaFold (green&cyan: crystal
structure / magenta: predicted model w/ residue_index modification).
<https://twitter.com/minkbaek/status/1417538291709071362> (2021).
- [22] Bryant, P. *et al. bioRxiv* 2021.09.15.460468 (2021).
- [23] Mosalaganti, S. *et al. bioRxiv* 2021.10.26.465776 (2021).

ACKNOWLEDGEMENTS

192 We thank Johannes Söding for providing computational re-
193 sources. Richard Evans, John Jumper, and Tim Green for
194 answering questions regarding AF2. Minkyung Baek for the
195 complex residue trick. Do-Yoon Kim for creating the Colab-
196 Fold logo. Enzo Guerrero-Araya and Jakub Kaczmarzyk for
197 providing bug fixes. Alon Markovich and Julia Varga for no-
198 tifying us about MSA quality issues. Harriet Alexander for
199 providing the TOPAZ proteins as a single file to download.
200 We thank all users for using ColabFold and reporting issues.

201 This work used the Scientific Compute Cluster at GWDG,
202 the joint data center of Max Planck Society for the
203 Advancement of Science (MPG) and University of Göttin-
204 gen. Milot Mirdita acknowledges the BMBF CompLifeSci
205 project horizontal4meta. Martin Steinegger acknowledges
206 support from the National Research Foundation of Ko-
207 rea grant [2019R1A6A1A10073437, 2020M3A9G7103933,
208 2021R1C1C102065, 2021M3A9I4021220]; New Faculty
209 Startup Fund and the Creative-Pioneering Researchers
210 Program through Seoul National University. Yoshitaka
211 Moriwaki acknowledges support from Platform Project for
212 Supporting Drug Discovery and Life Science Research (Basis
213 for Supporting Innovative Drug Discovery and Life Science
214 Research (BINDS)) from AMED under Grant Number
215 JP21am0101107. For this project, Sergey Ovchinnikov was
216 supported by the National Science Foundation under Grant
217 No. MCB2032259. Any opinions, findings, and conclusions
218 or recommendations expressed in this material are those of
219 the author(s) and do not necessarily reflect the views of the
220 National Science Foundation.

AUTHOR CONTRIBUTION

221 M.M., K.S., S.O. and M.S. performed research and program-
222 ming, M.M., S.O. and M.S. jointly designed the research and
223 wrote the manuscript. Y.M. provided the initial methodology
224 for hetero-complex modeling and created an installer for use
225 on local servers. L.H. provided initial benchmarking.

COMPETING INTERESTS

226 The authors declare no competing interests.

MATERIALS AND METHODS

Executing ColabFold

ColabFold is available as a set of Jupyter notebooks, to use on Google Colaboratory or users' local machines, as well as an easily installable command line application.

ColabFold notebooks ColabFold has four main Jupyter notebooks [24]: `AlphaFold2_mmseqs2` for basic use that supports protein structure prediction using (1) MSAs generated by MMseqs2, (2) custom MSA upload, (3) using template information, (4) relaxing the predicted structures using amber force fields [25], and (5) complex prediction. `AlphaFold2_advanced` for advanced users additionally supports (6) MSA generation using HMMer (same as AlphaFold-Colab), (7) the sampling of diverse structures by iterating through a series of random seeds (`num_samples`), and (8) control of AlphaFold2 model internals, such as changing the number of recycles (`max_recycle`), number of ensembles (`num_ensemble`), and enabling the stochastic part of the models via the (`is_training`) option. The latter enables dropout during inference, allowing the user to sample solutions from the uncertainty of the model [26] or the ambiguity of co-evolution constraints derived from the input MSA. `AlphaFold2_batch` for batch prediction of multiple sequences or MSAs. The batch notebook saves time by avoiding recompilation of the AlphaFold2 models ("Avoid recompiling during batch computation") for each individual input sequence. `RoseTTAFold` for basic use of RoseTTAFold that supports protein structure prediction using (1) MSAs generated by MMseqs2, (2) custom MSAs and (4) sidechain prediction using SCWRL4 [27]. The RoseTTAFold notebook also has an option use a slower but more accurate PyRosetta [28] folding protocol for structure prediction, using constraints predicted by RoseTTAFold's neural network.

ColabFold command line interface We initially focused on making ColabFold as widely available as possible through our Notebooks running in Google Colaboratory. To meet the demand for a version that runs on local users' machines, we released "LocalColabFold". LocalColabFold can take command line arguments to specify an input FASTA file, an output directory, and various options to tweak structure predictions. LocalColabFold runs on wide range of operating systems, such as Windows 10 or later (using Windows Subsystem for Linux 2), macOS, and Linux. The structure inference and energy minimization are accelerated if a CUDA 11.1 or later compatible GPU is present. LocalColabFold is available as free open-source software at github.com/YoshitakaMo/localcolabfold.

Recognizing the limitations of Google Colaboratory, we provide the `colabfold_batch` command line tool through the `colabfold` python package. This allows computing of tasks too large for Google Colab on users' own computer, e.g. predicting an entire proteome (Methods "Proteome benchmark"). It can be installed with `pip install colabfold`, followed by `pip install -U "jax[cuda]" -f https://storage.googleapis.com/jax-releases/jax_releases.html`. It can be used as `colabfold_batch input_file_or_directory output_directory`, supporting FASTA, A3M and CSV files as input.

Replacing MSA generation in AlphaFold2/RoseTTAFold with MMseqs2

Generating multiple sequence alignments for AlphaFold2 and RoseTTAFold is a time-consuming task. To improve their runtime, while maintaining a high prediction accuracy, we implemented optimized workflows using MMseqs2.

MSA generation by MMseqs2 ColabFold sends the query sequence to a MMseqs2 server [11]. It searches the sequence(s) with three iterations against the consensus sequences of the UniRef30, a clustered version of the UniRef100 [29]. We accept hits with an E-value of lower than 0.1. For each hit, we realign its respective UniRef100 cluster member using the profile generated by the last iterative search, filter them (Methods "New diversity aware filter") and add these to the MSA. This expanding search results in a speed up of ~10x as only 29.3 million cluster consensus sequence are searched instead of all 277.5 million UniRef100 sequences. Additionally, it has the advantages to be more sensitive since the cluster consensus sequences are used. We use the UniRef30 sequence-profile to perform an iterative search against the BFD/MGnify or ColabFoldDB using the same parameters, filters and expansion strategy.

New diversity aware filter To limit the number of hits in the final MSA we use the HHblits diversity filtering algorithm [8] implemented in MMseqs2 in multiple stages: (1) During UniRef cluster expansion, we filter each individual UniRef30 cluster before adding the cluster members to the MSA, such that no cluster-pair has a higher maximum sequence identity than 95% (`--max-seq-id 0.95`). (2) After realignment enable only the `--qsc 0.8` threshold and disable all other thresholds (`--qid 0 --diff 0 --max-seq-id 1.0`). Additionally, the qsc filtering is only used if least 100 hits were found (`--filter-min-enable 100`). (3) During MSA construction we filter again with the following parameters: `--filter-min-enable 1000 --diff 3000 --qid 0.0,0.2,0.4,0.6,0.8,1.0 --qsc 0 --max-seq-id 0.95`.

Here, we extended the HHblits filtering algorithm to filter within a given sequence identity bucket, such that it cannot eliminate redundancy across filter buckets. Our filter keeps the 3000 most diverse sequences in the identity buckets]0.0-0.2],]0.2-0.4],]0.4-0.6],]0.6-0.8] and]0.8-1.0]. In buckets containing less than 1000 hits we disable the filtering.

New MMseqs2 pre-computed index to support expanding cluster members MMseqs2 was initially built to perform fast many-against-many sequence searches. Mirdita *et al.* [11] improved it to also support fast single-against-many searches. This type of search requires the database to be index and stored in memory. `mmseqs createindex` indexes the sequences and stores all time-consuming-to-compute data structures used for MMseqs2 searches to disk. We load the index into the operating systems cache using `vmtouch` (github.com/hoytech/vmtouch) to allow calls to the different MMseqs2 modules to become near-overhead free. We extended the index to store, in addition to the already present cluster consensus sequences, all member sequences and the pairwise alignments of the cluster representatives to the cluster members. With these resident in cache, we eliminate the overhead of the remaining module calls.

ColabFold databases

342
343 AlphaFold2 requires over 2 terabyte of storage space for its
344 databases, which is a significant hurdle for many researchers.
345 We optimized its databases and additionally created another
346 large environmental sequence database.

347 **Reducing size of BFD/MGnify** To keep all required se-
348 quences and data structures in memory we needed to reduce
349 the size of the environmental databases BFD and MGnify, as
350 both databases together would have required ~517 GB RAM
351 for headers and sequences alone.

352 BFD is a clustered protein database consisting of ~2.2
353 billion proteins organized in 64 million clusters. MGnify
354 (2019_05) contains ~300 million environmental proteins. We
355 merged both databases by searching the MGnify sequences
356 against the BFD cluster representative sequences using MM-
357 seqs2. Each MGnify sequence with a sequence identity of
358 >30% and a local alignment that covers at least 90% of its
359 length is assigned to the respective BFD cluster. All unas-
360 signed sequences are clustered at 30% sequence identity and
361 90% coverage (`--min-seq-id 0.3 -c 0.3 --cov-mode 1 -s`
362 `3`) and merged with the BFD clusters, resulting in 182 million
363 clusters. In order to reduce the size of the database we fil-
364 tered each cluster keeping only the 10 most diverse sequences
365 using (`mmseqs filterresult --diff 10`). This reduced the
366 total number of sequences from 2.5 billion to 513 million, thus
367 requiring only 84 GB RAM for headers and sequences.

368 **ColabFoldDB** We built ColabFoldDB by expanding the
369 BFD/MGnify with metagenomic sequences from various en-
370 vironments. To update the database, we searched the pro-
371 teins from the SMAG (eukaryotes) [14], MetaEuk (eukary-
372 otes) [13], TOPAZ (eukaryotes) [15], MGV (DNA viruses) [16],
373 GPD (bacteriophages) [17] and updated version of MetaClust
374 [18] against the BFD/MGnify centroids using MMseqs2 and
375 assigned each sequence to the respective cluster if they have
376 a 30% sequence identity at a 90% sequence overlap (`-c 0.9`
377 `--cov-mode 1 --min-seq-id 0.3`). All remaining sequences
378 were clustered using `MMseqs2 cluster -c 0.9 --cov-mode`
379 `1 --min-seq-id 0.3` and appended to the database. We re-
380 move redundancy per cluster by keeping the most 10 diverse
381 sequences using (`mmseqs filterresult --diff 10`). The fi-
382 nal database consists of 209,335,865 million representative se-
383 quences and 738,695,580 members. See “Data availability” for
384 input files. We provide the MMseqs2 search workflow used in
385 the server (“MSA generation by MMseqs2”) as a standalone
386 script `colabfold_search.sh`.

387 **Template information** AlphaFold2 searches with HHsearch
388 through a clustered version of the PDB (PDB70 [8]) to find
389 the 20 top ranked templates. In order to save time, we use
390 MMseqs2 [10] to search against the PDB70 cluster represen-
391 tatives as a prefiltering step to find candidate templates. This
392 search is also done as part of the MMseqs2 API call on our
393 server. Only the top 20 target templates according to E-value
394 are then aligned by HHsearch. The accepted templates are
395 given to AlphaFold2 as input features. This alignment step is
396 done in the ColabFold client and therefore requires the subset
397 of the PDB70 containing the respective HMMs. The PDB70
398 subset and the PDB mmcif files are fetched from our server.
399 For benchmarking, no templates are given to ColabFold.

Modeling protein complexes with ColabFold

400 ColabFold offers protein complex folding through the spe-
401 cialized AlphaFold-multimer model and through residue-index
402 manipulation [3]. Here, we show the steps we took for Colab-
403 Fold to produce accurate protein complex predictions.

404 **Modeling of protein-protein complexes** We implemented
405 two protein complex prediction modes in ColabFold. One
406 based on AlphaFold-multimer [4] and one based on the residue
407 index manipulation of the original AlphaFold2 model. Baek
408 *et al.* [3] show that RoseTTAFold is able to model complexes,
409 despite being trained only on single chains. This is done by
410 providing a paired alignment and modifying the residue in-
411 dex. The residue index is used as an input to the models to
412 compute positional embeddings. In AlphaFold2, we find the
413 same to be true, although surprisingly the paired alignment
414 is often not needed (**Fig. 2c**). AlphaFold2 uses relative posi-
415 tional encoding with a cap at $|i-j| \geq 32$. Meaning, any pair
416 of residues separated by 32 or more are given the same relative
417 positional encoding. By offsetting the residue index between
418 two proteins to be > 32, AlphaFold2 treats them as separate
419 poly-peptide chains. ColabFold integrates this for modeling
420 complexes.

421
422 For homo-oligomeric complexes (**Fig. 3a**), the MSA is
423 copied multiple times for each component. Interestingly, it
424 was found that providing a separate MSA copy (padding by
425 gap characters to extend to other copies) to work significantly
426 better than concatenating left-to-right.

427 For hetero-oligomeric complexes (**Fig. 3b**), a separate MSA
428 is generated for each component. The MSA is paired according
429 to the chosen `pair_mode` (“MSA pairing for complex predic-
430 tion”). Since pLDDT is only useful for assessing local struc-
431 ture confidence, we use the fine-tuned model parameters to
432 return the PAE for each prediction. As illustrated in **Sup-**
433 **plementary Fig. 4**, the inter-PAE (predicted aligned error),
434 the predicted TM-score or interface TM-score (both derived
435 from PAE) can be used to rank and assess the confidence of
436 the predicted protein-protein interaction.

437 **MSA pairing for complex prediction** A paired MSA helps
438 AlphaFold2 to predict complexes more accurately only if or-
439 thologous genes are paired with each other. We followed a
440 similar strategy as Bryant *et al.* [22] to pair sequences accord-
441 ing to their taxonomic identifier. For the pairing we search
442 each distinct sequence of a complex against the UniRef100
443 using the same procedure as described in “MSA generation”.
444 We return only hits that cover all complex proteins within one
445 species and pair only the best hit (smallest e-value) with an
446 alignment that covers the query to at least 50%. The pairing
447 is implemented in the new MMseqs2 module `pairaln`.

448 For prokaryotic protein prediction, we additionally imple-
449 mented the protocol described in [3] to pair sequences based
450 on their distances in the genome as predicted from the UniProt
451 accession numbers.

452 **Taxonomic labels for MSA pairing** To pair MSAs for com-
453 plex prediction, we retrieve for each found UniRef100 member
454 sequence the taxonomic identifier from the NCBI taxonomy
455 [30]. The taxonomic labels are extracted from the lowest com-
456 mon ancestor field (“common taxon ID”) of each UniRef100
457 sequence from the `uniref100.xml` (2021_03) file.

Speeding up AlphaFold2's model evaluation

Our efforts in speeding up AlphaFold2's MSA generation yielded large improvements in its runtime. However, we discovered multiple opportunities within AlphaFold2 to speed up its model inference, without sacrificing (or only sacrificing very little) of its accuracy.

Avoid recompiling AlphaFold2 models The AlphaFold2 models are compiled using JAX [31] to optimize the model for specific MSA or template input sizes. When no templates are provided, we compile once and, during inference, replace the weights from the other models, using the configuration of model 5. This saves 7 minutes of compile time. When templates are enabled, model 1 is compiled and weights from model 2 are used, model 3 is compiled and weights from models 4 and 5 are used. This saves 5 minutes of compile time. If the user changes the sequence or settings, without changing the length or number of sequences in the MSA, the compiled models are reused without triggering recompilation.

Avoid recompiling during batch computation In order to avoid AlphaFold2 model recompilation for every protein AlphaFold2 provides a function to add padding to the input MSA and templates called `make_fixed_size`. However, this is not exposed in AlphaFold2. We used the function in our batch notebook as well as in our command line tool `colabfold_batch`, in order to maximize GPU utilization and minimize the need of model recompilation. We sort the input queries by sequence length and process them in ascending order. We pad the input features by 10% (by default). All sequences that lie within the query length and an additional 10% margin do not require to be recompiled, resulting in a large speed up for short proteins.

Recycle count AlphaFold2 improves the predicted protein structure by recycling (by default) 3 times, meaning the prediction is fed multiple times through the model. We exposed the recycle count as a customizable parameter as additional recycles can often improve a model (**Supplementary Fig. 6**) at the cost of a longer runtime. We also implemented an option to specify a tolerance threshold to stop early. For some designed proteins without known homologous sequences, this helped to fold the final protein (**Supplementary Fig. 5**).

Speed-up of predictions through early stop AlphaFold2 computes five models through multiple recycles. We noted that for prediction of high certainty (> 85 pLDDT), all five models would often produce structures of very similar confidence, for some even without or with less than 3 of recycles. In order to speed up the computation we added a parameter to define an early stop criterion that halts additional model inferences and stops recycling if a given pLDDT or (interface) pTMscore threshold is reached.

Exposing advanced features

In our investigation of AlphaFold2's internals, we realized that we could expose many knobs that might be usefully to researchers trying to explore AlphaFold2's full potential.

Sampling of diverse structures To reduce memory requirements, only a subset of the MSA is used as input to the model. AlphaFold2, depending on model configuration, subsamples the MSA to a maximum of 512 cluster centers and 1024 "extra" sequences. Changing the random seed can result in different

cluster centers and thus different structure predictions. ColabFold provides an option to iterate through a series of random seeds, resulting in structure diversity. Further structure diversity can be generated by using the original or fine-tuned (`use_ptm`) model parameters and/or enabling (`is_training`) to activate the stochastic (dropout) part of model. Enabling the latter, can be used to sample an ensemble of models for the uncertain parts of the structure prediction.

Custom MSAs ColabFold allows researchers to upload their own MSAs. Any kind of alignment tool can be used to generate the MSA. The uploaded MSA can be provided in aligned FASTA, A3M, STOCKHOLM or Clustal format. We convert the respective MSA format into A3M format using the `reformat.pl` script from the HH-suite [8].

Lightweight 2D structure renderer For visualization, we developed a matplotlib [32] compatible module for displaying the 3D ribbon diagram of a protein structure or complex. The ribbon can be colored by residue index (N to C terminus) or by a predicted confidence metric (such as pLDDT). For complexes, each protein can be colored by chain ID. Instead of using a 3D renderer, we instead use a 2D line plotting based technique. The lines that make up the ribbon are plotted in the order in which they appear along the z-axis. Furthermore, we add shade to the lines according to the z-axis. This creates the illusion of a 3D rendered graphic. The advantage over a 3D renderer is that the images are very lightweight, can be used in animations and saved as vector graphics for lossless inclusion in documents. As the 2D renderer is not interactive, we additionally included a 3D visualization using py3Dmol [33] in the ColabFold notebooks.

Benchmarking ColabFold

We show with multiple datasets that ColabFold does not sacrifice accuracy for its much faster runtimes.

Benchmark with CASP14 targets We compared AlphaFold-Colab and AlphaFold2 (commit `b88f8da`) against ColabFold using all CASP14 [2] targets. ColabFold-AlphaFold2 (commit `2b49880`) used UniRef30 (2021_03) [34] and the BFD/MGnify or ColabFoldDB. ColabFold-RoseTTAFold (commit `ae2b519`) was executed with papermill (`github.com/nteract/papermill`) using the PyRosetta protocol [28]. ColabFold-RoseTTAFold-BFD/MGnify and ColabFold-AlphaFold2-BFD/MGnify used the same MSAs. AlphaFold-Colab used the UniRef90 (2021_03), MGnify (2019_05) and the small BFD. AlphaFold2 used the `full_dbs` preset with and default databases downloaded with the `download_all_data.sh` script. The 65 targets contain 91 domains, among these are 20 FM-targets with 28 domains. We compared the predictions against the experimental structures using TMalign [35].

Measuring run-times for CASP14 benchmark To provide more accurate run times we split MSA generation and model inference measurements. MSA generation times were repeated five times and averaged.

ColabFold was executed using `colabfold_batch`. The MM-seqs2 server which computes MSAs for ColabFold has 2x14 core Intel E5-2680v4 CPUs and 768 GB RAM. Each generated MSA was processed by a single CPU-core. Runtimes were computed from server logs.

573 AlphaFold2 MSA generation runtimes were measured by
574 running AlphaFold2 without models (providing an empty
575 string to the `--model_names` parameter) on the same 2x14
576 core Intel E5-2680v4 CPUs and 768 GB RAM system. The
577 AlphaFold2 databases were stored on a software-RAID5 com-
578 posed of six Samsung 970 EVO Plus 1TB NVMe drives. Run-
579 times for AlphaFold2 were taken from the `features` entry of
580 the `timings.json` file. For a fair comparison, AlphaFold2 was
581 modified to allow HMMer and HHblits to access one CPU core.

582 All ColabFold and AlphaFold2 model inference runtime
583 measurements were done on systems with 2x16 core Intel
584 Gold 6242 CPUs with 192 GB RAM and 4x Nvidia Quadro
585 RTX5000 GPUs. Only one GPU was used in each run.

586 ColabFold-RoseTTAFold-BFD/MGnify and ColabFold-
587 AlphaFold2-BFD/MGnify used the same MSAs, runtimes are
588 shown only once.

589 AlphaFold-Colab was executed in the browser using a
590 Google Colab Pro account. Times for homology search were
591 taken from the notebook output cell “Search against genetic
592 databases” cell. The JackHMMer search uses 8 threads.

593 **Complex benchmark** We compare predictions of seventeen
594 ClusPro [4, 12] targets to their native structures using DockQ
595 [36]. We used `colabfold_batch` (commit 45ad0e9) with
596 BFD/MGnify in residue-index manipulation- and AlphaFold-
597 multimer mode to predict structures. We use MSA pairing as
598 described in “MSA pairing for complex prediction” and also
599 add unpaired sequences. Models are ranked by predicted in-
600 terface pTMScore as returned by AlphaFold-multimer. The
601 DockQ AlphaFold-multimer reference numbers were provided
602 by Richard Evans.

603 **Proteome benchmark** We predict the proteome of *M. jan-*
604 *naschii*. Of the 1787 proteins we exclude the 25 proteins longer
605 than 1000 residues, leaving 1762 proteins of 268 aa average
606 length. With the `colabfold_search` wrapper to MMseqs2
607 we search against the ColabFoldDB (“ColabFoldDB”) in 113
608 min on a system with an AMD EPYC 7402P 24-core CPU (no
609 hyperthreading) and 512GB RAM. MMseqs2 had a maximum
610 resident set size of 308 GB during the search. We then predict
611 the structures on a single Nvidia Titan RTX with 24 GB RAM
612 in 46 h using only MSAs (no templates). For each query we
613 stop early if any recycle iteration reaches a pLDDT of at least
614 85. Early stopping results in a speed-up of 3.7× over default
615 and 4.8× over always recompiling. AlphaFold2 (`reduced_dbs`)
616 was ran with the `reduced_dbs` preset and no template infor-
617 mation was used. We changed the AlphaFold2 source code to
618 utilize all CPU cores during the homology search.

619 AlphaFold2 (`reduced_dbs`, v2.1.1), ColabFold (commit
620 f5d0cec) default and ColabFold Stop ≥ 85 have an average
621 pLDDT of 90.68, 90.22 and 89.33 respectively for 50 ran-
622 domly sampled proteins. These are the same proteins that
623 were used to extrapolate the run-time of AlphaFold2. Over
624 all predictions, the pLDDTs for the *M. jannaschii* proteome
625 downloaded from the AlphaFoldDB, ColabFold default and
626 ColabFold Stop ≥ 85 are 89.75, 89.38 and 88.77, respectively.

CODE AVAILABILITY

627 ColabFold is free open-source software (MIT) and avail-
628 able at github.com/sokrypton/ColabFold. A locally in-
629 stallable version is available at [github.com/YoshitakaMo/](https://github.com/YoshitakaMo/localcolabfold)
630 `localcolabfold`. The ColabFold development version shown
631 in this manuscript is available at [github.com/konstin/](https://github.com/konstin/ColabFold)
632 `ColabFold`. The ColabFold server components are free
633 open-source software (GPLv3) and available at [github.com/](https://github.com/soedinglab/mmseqs2-app)
634 `soedinglab/mmseqs2-app`. MMseqs2 is free open-source soft-
635 ware (GPLv3) and available at mmseqs.com.

DATA AVAILABILITY

636 ColabFold databases are freely (CC-BY-SA 4.0) available at
637 colabfold.mmseqs.com.
638 MSAs and structures produced during benchmarking:
639 wwwuser.gwdg.de/~compbiol/colabfold/manuscript
640 Input databases used for building ColabFold databases:
641 UniRef30: uniclust.mmseqs.com
642 BFD: bfd.mmseqs.com
643 MGnify: [ftp.ebi.ac.uk/pub/databases/metagenomics/](https://ftp.ebi.ac.uk/pub/databases/metagenomics/peptide_database/2019_05)
644 `peptide_database/2019_05`
645 PDB70: [wwwuser.gwdg.de/~compbiol/data/hhsuite/](https://wwwuser.gwdg.de/~compbiol/data/hhsuite/databases/hhsuite_dbs)
646 `databases/hhsuite_dbs`
647 MetaEuk: wwwuser.gwdg.de/~compbiol/metaeuk/2019_11/
648 `MetaEuk_preds_Tara_vs_euk_profiles_uniqs.fas.gz`
649 SMAG: [www.genoscope.cns.fr/tara/localdata/data/](https://www.genoscope.cns.fr/tara/localdata/data/SMAGs-v1/SMAGs_v1_concat.faa.tar.gz)
650 `SMAGs-v1/SMAGs_v1_concat.faa.tar.gz`
651 TOPAZ: osf.io/gm564
652 MGv: [portal.nersc.gov/MGV/MGV_v1.0_2021_07_08/mgv_](https://portal.nersc.gov/MGV/MGV_v1.0_2021_07_08/mgv_proteins.faa)
653 `proteins.faa`
654 GPD: [ftp.ebi.ac.uk/pub/databases/metagenomics/](https://ftp.ebi.ac.uk/pub/databases/metagenomics/genome_sets/gut_phage_database/GPD_proteome.faa)
655 `genome_sets/gut_phage_database/GPD_proteome.faa`
656 Further datasets used for benchmarking ColabFold:
657 PFAM (Pfam-A.seed.gz & Pfam-A.full.gz):
658 ftp.ebi.ac.uk/pub/databases/Pfam/releases/Pfam34.0
659 *M. jannaschii* proteome:
660 uniprot.org/proteomes/UP000000805
661 [ftp.ebi.ac.uk/pub/databases/alphafold/v1/](https://ftp.ebi.ac.uk/pub/databases/alphafold/v1/UP000000805_243232_METJA_v1.tar)
662 `UP000000805_243232_METJA_v1.tar`

REFERENCES

- [24] Kluyver, T. *et al.* Jupyter Notebooks - a publishing format for reproducible computational workflows. In *Positioning and Power in Academic Publishing: Players, Agents and Agendas*, 87–90 (IOS Press, 2016).
- [25] Eastman, P. *et al.* *PLoS Comput. Biol.* **13**, 1–17 (2017).
- [26] Gal, Y. & Ghahramani, Z. *arXiv* 1506.02142 (2016).
- [27] Krivov, G. G. *et al.* *Proteins* **77**, 778–795 (2009).
- [28] Chaudhury, S. *et al.* *Bioinformatics* **26**, 689–691 (2010).
- [29] Suzek, B. E. *et al.* *Bioinformatics* **31**, 926–932 (2015).
- [30] Federhen, S. *Nucleic Acids Res.* **40**, D136–D143 (2012).
- [31] Bradbury, J. *et al.* JAX: composable transformations of Python+NumPy programs (2018).
- [32] Hunter, J. D. *Comput. Sci. Eng.* **9**, 90–95 (2007).
- [33] Rego, N. & Koes, D. *Bioinformatics* **31**, 1322–1324 (2015).
- [34] Mirdita, M. *et al.* *Nucleic Acids Res.* **45**, D170–D176 (2017).
- [35] Zhang, Y. & Skolnick, J. *Nucleic Acids Res.* **33**, 2302–2309 (2005).
- [36] Basu, S. & Wallner, B. *PLoS One* **11**, e0161879 (2016).