

CryoFold: determining protein structures and ensembles from cryo-EM data

Mrinal Shekhar¹, Genki Terashi², Chitrak Gupta^{3,4}, Daipayan Sarkar^{2,3}, Gaspard Debussche⁵, Nicholas J. Sisco^{3,6}, Jonathan Nguyen^{3,4}, Arup Mondal⁹, James Zook^{3,4}, John Vant^{3,4}, Petra Fromme^{3,4}, Wade D. Van Horn^{3,6}, Emad Tajkhorshid¹, Daisuke Kihara (c)^{2,8}, Ken Dill (c)⁷, Alberto Perez (c)⁹, and Abhishek Singharoy (c)^{3,4}

¹Center for Biophysics and Quantitative Biology, Department of Biochemistry, NIH Center for Macromolecular Modeling and Bioinformatics, Beckman Institute for Advanced Science and Technology, University of Illinois at Urbana-Champaign, Urbana, Illinois, 61801, USA.; ²Department of Biological Sciences, Purdue University, West Lafayette, IN 47907, USA.; ³The School of Molecular Sciences, Arizona State University, Tempe, AZ 85287, USA.; ⁴Biodesign Institute Center for Structural Discovery, Arizona State University, Tempe, AZ 85281, USA.; ⁵Department of Mathematics and Computer Sciences, Grenoble INP, 38000 Grenoble, France.; ⁶The Biodesign Institute Virginia G. Piper Center for Personalized Diagnostics, Arizona State University, Tempe, AZ 85281, USA.; ⁷Laufer Center for Physical and Quantitative Biology, Stony Brook University, Stony Brook, New York 11794, United States; ⁸Department of Computer Science, Purdue University, West Lafayette, IN 47907, USA.; ⁹Chemistry Department, University of Florida, Gainesville, Florida, 32611, USA

Cryo-EM is a powerful method for determining protein structures. But it requires computational assistance. Physics-based computations have the power to give low-free-energy structures and ensembles of populations, but have been computationally limited to only small soluble proteins. Here, we introduce CryoFold. By integrating data of varying sparsity from electron density maps of 3–5 Å resolution with coarse-grained physical knowledge of secondary and tertiary interactions, CryoFold determines ensembles of protein structures directly from sequence. We give six examples showing its broad capabilities, over proteins ranging from 72 to 2000 residues, including membrane and multi-domain proteins, and including results from two EMDB competitions. The ensembles CryoFold predicts starting from the density data of a single known protein conformation encompass multiple low-energy conformations, all of which are experimentally validated and biologically relevant.

Protein folding | CryoEM | free-energy | Molecular dynamics flexible fitting | MAINMAST | MELD | Ensemble refinement

Cryo-electron microscopy (cryo-EM) is a powerful tool for determining the structures of biomolecules. It serves a niche – such as large complexes or membrane proteins or molecules that are not easily crystallizable – that traditional methods, such as X-ray diffraction, electron or neutron scattering, or NMR often cannot handle. Routine cryo-EM structure determination has a number of components: the experiment produces raw data in the form of single-particle images, correction and processing of this data recovers an electron density map, and finally molecular modeling is required to determine structures from the map. Currently, there are two broad classes of methods for molecular modeling. First, established algorithms for refining X-ray crystallography or NMR structures, such as *Phenix.realspacerefine* or REFMAC (1), are often used, even for ensemble determination (2), but offer complete models with the highest-resolution density data. Cryo-EM studies commonly produce lower-resolution data. Second, *integrative approaches* that leverage data from multiple types of experiments (3) to find structures compatible with the data. The challenge here is that cryo-EM data is often *heterogeneous*, meaning that some parts of a protein structure are well-determined by the data while others are more poorly defined.

For computational modeling, the changing resolution poses the need for extensive conformational sampling and the need to identify which conformations amongst all that fit the lower resolution regions are most biophysically relevant. The size of

the search space is large and grows non-linearly with system size (4). Physics-based modeling, such as molecular dynamics (MD) simulations, can give proper thermodynamic weights for choosing among the different conformational populations. But, we need efficient ways of sampling using physics based approaches. Most MD is used for exploring dynamics around an experimental structure and for automated model refinement (5, 6). Yet, large conformational changes, such as those relevant in many biological processes, remain inaccessible to MD(7–9) – it is computationally expensive. In structure determination, the end structure is unknown, so collective variables to accelerate the process are not an option (10). Therefore, MD is augmented with external information such as evolutionary covariance (11, 12) and homology-based starting models (13, 14), or with advanced sampling methods based on Bayesian inference (15–17) and specialized hardware (18), which improve the speed of structure prediction by 10 to 100-fold over brute-force simulations. Notwithstanding this improvement, the prediction of protein fragments beyond 115 residues remain a bottleneck for physics-based methods (19). Fragment search and fitting schemes are successful in resolving the EM map(20), but they require at least 70% of the C α atoms placed correctly (21–23), and for membrane systems, such refinements also leverage MD simulations (24). However, the bioinformatic augmentations to MD introduce new discrepancies that are often refractory to automated fixes (23, 25), warranting our developments.

Here, we describe CryoFold, an integrative atomistic-physical algorithm that derives ensemble of folded protein structures from cryo-EM data. Illustrated in (Fig. 1), CryoFold is a combination of three methods: (1) MAINMAST(26), MAINchain Model trAcing from Spanning Tree – a method that generates the trace of the connected peptide chain when provided with EM data, (2) ReMDFF(27), Resolution exchange Molecular Dynamics Flexible Fitting – a MD method for refining protein conformations from electron-density maps, and (3) MELD(15, 28), Modeling Employing Limited Data – a Bayesian folding and refolding engine that can work from insufficient data to accelerate the MD sampling of rare events such as those needed for protein folding. The guidance from experimental data allows MD simulations to fold models with well beyond 115 residues, including transmembrane systems

¹M.S.(Author One) and G.T. (Author Two) contributed equally to this work.

70 and asymmetric multi-protein complexes. More importantly, 130
71 the free energy description of folded and unfolded popula- 131
72 tions accessible to MELD enables the exhaustive sampling 132
73 of structures that are representative of different metastable 133
74 states. Thus, starting with the structural data from a par- 134
75 ticular protein conformation, CryoFold predicts on one hand, 135
76 the energetically favorable ensemble of structures that are 136
77 consistent with the data, while on the other hand, discovers 137
78 multiple new low-energy protein states in the vicinity of the fit-
79 ted model. Going beyond the determination of one stationary
80 structure, CryoFold offers the opportunity to combine all the
81 predicted structures into a model conformational transition
82 pathway, where the new states are also validated and re-refined
83 against orthogonal NMR, X-ray crystallography or cryo-EM
84 data.

85 Starting with density maps of resolution 5.0 Å and higher,
86 first, MAINMAST is employed to derive a chain trace of C_α
87 atoms. Then we use this trace as a template to iterate between
88 MELD and ReMDFF. While MELD explores a large conforma-
89 tional space, visiting multiple plausible secondary structures
90 consistent with the MAINMAST template, ReMDFF simula-
91 tions refine the protein backbone and sidechain conformations
92 to fit to the density map for each one of the assumed secondary
93 structures (6). Taken alone, ReMDFF fits models into electron
94 density features, but fails to explore the variations in secondary
95 structures(27). MELD addresses this issue by partial folding,
96 unfolding and reformation of secondary structures(15, 28),
97 using the coarse physical information (CPI) available on web-
98 servers(15, 29); for example, based on their sequences, proteins
99 prefer specific fractions of hydrophobic interactions, β -strand
100 pairing and secondary structures to minimize frustration (Fig.
101 2A). Consequently, a hybrid iterative MELD-ReMDFF ap-
102 proach allows the determination of complete all-atom models
103 from sequence information merged with available structural
104 data of varying coarseness. For intermediate to low-resolution
105 data (lesser than 5 Å) wherein C-alpha tracing is unreliable
106 (30), the MAINMAST step can be avoided. Nonetheless, if
107 successful, the search template derived from backbone tracing
108 almost always accelerates convergence of CryoFold.

109 We report data-guided structural ensembles for six different
110 examples here, for proteins from 72 to 618 residues, extending
111 to multi-protein complexes of up to 2000 residues, and across
112 both soluble and membrane systems. CryoFold overcomes the
113 sampling limitations of traditional MD predictions, producing
114 high-quality structural models: it offers a high radius of con-
115 vergence in the range of 50 Å, refining soluble and transmem-
116 brane structures with consistently > 90% favored backbone
117 and sidechain statistics, and high EMRinger scores (31). The
118 results are independent of the initial estimated conformation
119 and consistent with physics and stereochemistry, highlighted
120 through results in 2016 and 2019 EMD competition. The
121 hybrid protocol is available through a python-based graphical
122 user interface with a video tutorial.

123 Results

124 We describe six systems, chosen to represent the different
125 bottlenecks in the three component methods of the CryoFold
126 pipeline. At any given resolution, the accuracy of CryoFold
127 predictions depends on: (1) quality of C_α traces by MAIN-
128 MAST, (2) variations in secondary structure within the MELD
129 ensemble, and (3) convergence of ReMDFF. Three are soluble

130 proteins, with varying degrees of local resolution in the density
131 maps. One is from the 2019 EMD competition challenge, 132
133 in which the data was provided at three different resolutions. 134
135 One was a large asymmetric multi-protein complex that al-
136 lowed us to test how big a structure we could handle. And,
137 one was a transmembrane system, to see if MELD's aqueous
implicit-solvent model would be sufficient for the membrane
environment.

A. Proof of principle on a small known protein. In this case, 138
we began with a synthetic map of ubiquitin, a small 72-residue 139
protein. Ubiquitin is a good test system because, on the one 140
hand, it is small enough to fold computationally, and yet on the 141
other hand its experimental folding time is in the millisecond 142
range, so it been hard to fold by brute force MD (32), and 143
even, to a lesser extent, by the MELD approach (15). From 144
the known X-ray crystal structure of ubiquitin, we generated a 145
synthetic electron density at 3.0 Å resolution (33), and asked 146
if CryoFold could correctly recover the X-ray structure. We 147
found that only two MELD-ReMDFF iterations (Fig. 2B) 148
were needed to give a model having an RMSD difference of 149
2.53 Å from the crystal structure (PDB id: 1UBQ, see Table 150
S1). 151

**B. Test on a soluble lipoprotein with a uniformly high-resolu- 152
tion data.** Francisella lipoprotein Flpp3 is a 108 amino acids 153
long membrane-interacting protein that serves as a target for 154
drug development against tularemia(34). In this case, we 155
had two datasets: one at high resolution (1.8 Å) from our 156
Serial Femtosecond X-ray (SFX) crystallography experiments 157
of Flpp3 (See Supplementary Information and (35), and a 158
synthetic one at low resolution (5.0 Å). The point of this 159
test was to see if we could use the low-resolution data to 160
achieve the high-resolution structure. For both sets, we used 161
MAINMAST (26) to introduce the C_α traces as constraints 162
for MELD (Fig. 3A,B). Convergent ensembles derived from 163
this MAINMAST-guided MELD step were then refined by 164
ReMDFF to improve the sidechains until the density was 165
resolved with models of reliable geometry. 166

We found that one iteration of the MELD-ReMDFF cycle 167
sufficed to resolve an all-atom model of Flpp3 from the 168
SFX density, with accurate sidechain conformations, secondary 169
and tertiary structure assignments (structural statistics sum- 170
marized in Table S2). At 5 Å resolution MAINMAST pro- 171
duced low quality backbone traces (Fig. 3B). Remarkably, 172
even these low quality C_α traces, were enough for MELD- 173
ReMDFF to successfully produced models comparable to our 174
high-resolution refinements. After two MELD-ReMDFF itera- 175
tions, the best structure obtained was within 2.29 Å RMSD 176
from the SFX model. The MELD-only predictions modelled 177
the β -sheets accurately, they failed to accurately converge 178
on all helices (SI Fig. S1). For example, a 4-turn helix was 179
underestimated to contain only 2-3 turns. However, guidance 180
by the density map in CryoFold recovered these turns in both 181
the high and low resolution cases (Tables S2 and S3). Thus, 182
the Flpp3 test shows that the CryoFold trio of methods gives 183
accurate structures for longer chains than is otherwise possible 184
with either one of these methods. 185

Here, we are also able to test an important aspect of physics- 186
based structure determination, namely whether we can gener- 187
ate proper conformational ensembles, not just single average 188
structures. The quality of the CryoFold ensembles is accessed 189

190 against a set of 20 NMR models of Flpp3 (34) by looking at the
191 conformation of key residues (Y83,K35 and D4) responsible
192 for binding tularemia drugs (Fig. S2A). Upon projecting the
193 ensemble of 50 lowest-energy CryoFold structures onto a space
194 defined by the distance between Y83-K35 & Y83-D4, where
195 closed Flpp3 is represented by (Y83-K35 <5.00 Å & Y83-D4
196 > 10.00 Å), and open Flpp3 implies (Y83-K35 >10.00 Å &
197 Y83-D4 < 5.00 Å), all the major conformational states seen
198 in the NMR experiments have been recovered (Fig. S2B).
199 Thus, extending beyond the prediction of a single stationary
200 structure, the cluster of low-energy conformations predicted
201 by CryoFold captures both the open and closed conformations,
202 starting only with data from the closed state. The classifica-
203 tion of structural ensembles based on projections onto the
204 distance space requires *a priori* knowledge of the structural
205 features of all the major states in the ensemble. In an alternate
206 scheme that does not require such knowledge, the models were
207 classified based on their Rosetta-energy and RMSD relative to
208 the crystal structure (36). Rosetta is chosen as a benchmark
209 due to its use of energy functions analogous to the CHARMM
210 or AMBER force fields (37, 38) in MELD and MDFF (39).
211 In this energy space, the ensemble of structures derived from
212 Rosetta-EM visited almost all the states of Flpp3 observed
213 in NMR, while CryoFold recovered only a minimum number
214 of these states at 1.8 Å resolution (Fig. S3). In contrast, for
215 the (5.00 Å) regime, CryoFold shows a markedly better perfor-
216 mance with predictions overlapping with the majority of NMR
217 intermediates, as well as consistently determining lower energy
218 structures than Rosetta-EM. Thus, extended sampling benefits
219 of CryoFold is apparent in fuzzier data sets. Here, a broader
220 segment of the protein folding funnel is accessed by MELD,
221 recovering models from the poor initial guesses generated by
222 MAINMAST (Fig. S4). Taken together, the ubiquitin and
223 Flpp3 examples establish CryoFold as an enhanced sampling
224 tool for resolving multiple metastable states of proteins with
225 > 100 residues, guided only by a single experimental data set
226 at 3-5 Å

227 **C. Test on soluble domains of a membrane protein with het-**
228 **erogeneous-resolution data.** We look at the cytoplasmic do-
229 main of a large trans-membrane protein, TRPV1, a heat-
230 sensing ion channel (592 amino acids long). The point of
231 this test is that the data is highly heterogeneous, with ex-
232 perimental electron densities ranging between 3.8 to 6.0 Å
233 (40, 41), as determined by Resmap (42). Furthermore, TRPV1
234 has two apo-structures deposited in the RCSB database, one
235 with moderately resolved transmembrane helices and cytoplasmic
236 domains (41) (pdb id:3J5P, EMDataBank: EMD-5778),
237 and another with highly-resolved transmembrane helices (pdb
238 id:5IRZ, EMDataBank: EMD-8118) but with the cytoplasmic
239 regions, particularly the β -sheets, less resolved than in 3J5P.

240 CryoFold was employed to regenerate these unresolved
241 segments of the cytoplasmic domain from the heterogeneous
242 lower-resolution data of 5IRZ. We compare the CryoFold model
243 to the reported 3J5P structure (Fig. 4), where these domains
244 are much better resolved showing clear patterns of β -strands.
245 The final model was observed to be at an RMSD of 3.41 Å
246 with a CC of 0.74 relative to 5IRZ. The same model with some
247 loops removed for consistency with the EMD-5778 density
248 produced an RMSD of 2.49 Å and CC of 0.73 with respect
249 to 3J5P. Taken together, models derived from the CryoFold
250 refinement of 5IRZ capture in atomistic details the highly

251 resolved features of this density, yet without compromising
252 with the mid-resolution cytoplasmic areas where it performs
253 as well as the 3J5P model (Table S4).

254 TRPV1 was part of the 2016 Cryo-EM modeling challenge
255 where only ReMDFF was used (43). Presented in Table S5, our
256 updated CryoFold model of TRPV1 (model no. 4), represents
257 the the top - 20% of the submissions with > 90% Ramachan-
258 dran favored statistics, and an EMRinger score of 2.54. This
259 model is vastly refined over the originally reported structure
260 with a score of 1.75, and our previous submission at 2.25. The
261 improvement is attributed solely to the higher-quality β -sheet
262 models that is now derived from the enhanced sampling ob-
263 tained by running MELD and ReMDFF in tandem. Starting
264 with a random coil as search model (Fig. 4B), the recovery of
265 these β -sheets is highly improbable with the limited conforma-
266 tional space that MDFF visits. Addressing this issue, MELD
267 invokes a multi-replica temperature exchange scheme, wherein
268 at high replica indices it samples many distinct structures that
269 have short lifetimes (44). At the lower-temperature replica
270 a stronger coupling with the data is achieved, and these struc-
271 tures are folded into a smaller number of long-lived clusters,
272 each with varying degrees of native contacts and secondary
273 structure (Fig.S5). Thus, unlike MDFF, MELD allows for a
274 search of structural motifs constrained by features in the data.
275 When these methods are combined within CryoFold, both the
276 backbone and sidechain geometries are refined to capture rare
277 secondary structural changes, enabling the determination of
278 TRPV1's labile β -sheets.

279 An analysis of the CryoFold ensembles reveal partial un-
280 folding of the *beta*-sheets in the soluble domains of TRPV1
281 with around 3-4% of the structures presenting incomplete
282 *beta*-sheets, akin to the model originally submitted with 3J5P
283 (Fig.S5C). Partial unfolding of these regions have not been
284 attributed to any functional implications in TRPV1,
285 though some peripheral evidence of functional advantages
286 from unfolding exist in TRPV3 channels (45). The β -sheets
287 and loops from the soluble domains form the inter-protomer
288 interface within the tetrameric channel. Secondary structural
289 changes at these interfaces, triggers coupling between cyto-
290 plasmic and transmembrane domains, priming the channel
291 for opening. Such changes, though rare, are indeed appar-
292 ent in our MELD assignments. Therefore, the ensemble of
293 structures and not merely a single model that CryoFold offers,
294 opens the door to analyzing a number of distinct folded and
295 unfolded conformations, all of which contribute to the same
296 density map (46-48). Also evident from the TRPV1 case
297 study, we can generate such atomistic ensembles with data
298 of low local-resolution, yet with accuracy commensurate to
299 structures derived from higher resolution density maps.

300 **D. Tests on apoferritin at three different resolutions from the**
301 **2019 EMDB modeling challenge.** The EMDB competition is a
302 community-wide effort to assess the limits of structure predic-
303 tion using cryo-EM data. Here we were tasked to determine the
304 structure of an apoferritin monomer using data at 1.8, 2.3 and
305 3.1 Å resolution. Following an initial tracing by MAINMAST
306 on the monomeric map, it took two iterations for CryoFold to
307 arrive at the final model for the first two resolutions, and three
308 iterations for the third map. In total 17 teams participated
309 in the 2019 competition that focused primarily on ab-initio
310 structure determination, and all the results are reported on
311 the EMDB website (49). CryoFold (team 73) models were

312 independently assessed to be high accuracy (Fig. S6 (scale
313 labeled in green)), specifically for three different categories of
314 scores: Reference-free, EM-map and target-structure scores.
315 The results were robust over the narrow range of resolutions
316 tested, earning us the top rank for multiple entries (48). Comparability
317 with respect to the target structures is almost always
318 very high, as also reflected in commensurately high Fourier
319 Shell Coefficient ($FSC = 0.5$) and cross correlations with the
320 experimental map. Another noticeable strength is the strong
321 EMRinger scores of the MD-based refinement, very similar
322 to MDFF's performance in the 2016 competition (43). A
323 relatively new measure to evaluate mainchain geometry and
324 to identify areas of probable secondary structure based on
325 C-Alpha geometry, called CaBLAM (50) also found the CryoFold
326 models to be favorable. One limitation however, is the
327 increased number of Ramachandran outliers observed in the
328 CryoFold and MDFF determined structures, which implicates
329 the assumptions of classical CHARMM-type force fields(43).
330 Our recently developed neural network potentials have already
331 been useful to circumvent this issue (43, 46).

332 **E. Test on a large multi-chain protein complex with mid-resolution data.**
333 A grand challenge for cryo-EM is to determine
334 structures of multi-chain complexes. Symmetry is used wherever
335 possible, e.g., in viruses or homo-oligomeric membrane
336 proteins (45, 51). However, most protein-protein or protein-
337 nucleic acid complexes are asymmetric. Our test here is
338 whether CryoFold could obtain the structure in an asymmetric
339 complex. We focused on ATP synthase. It contains 31
340 chains. Recently Murphy et al. reported 30 distinct conformations
341 of this motor at 2.7-4.3 Å resolution (52). Similar to the Flpp3
342 and TRPV1 cases, here the ensemble computed by CryoFold
343 correctly captured the low-lying states of the multi-chain system
344 in addition to the target 6RET conformation. For simplicity,
345 we have removed the transmembrane *c*-ring of this system;
346 the transmembrane challenge will be addressed in the next section.
347

348 Seven of the reported thirty models by Murphy et al. included
349 overall deformations of the system without rotation of the *c*-ring.
350 Using RMSD matrices (Fig. S7A), these structures were clustered
351 in 4 distinct states (States I: 6RET; II: 6RDQ, 6RDR; III: 6RDK,
352 6RDL; and IV: 6RDW, 6RDX). Remarkably, all these four states
353 are identifiable in an RMSD matrix of 220 MELD structures
354 within CryoFold (Fig. 5B). States II, III and IV from MELD
355 are initially at RMSD 7.6, 12.0 and 8.4 Å from 6RET respectively
356 (Fig: S7B). After MDFF refinements, structures are consistent
357 with experimental models from Murphy et al. listed for states
358 II, III and IV were refined to RMSD values of 2.1, 2.8, and 1.8 Å
359 relative to the target models (Fig. 5C, S7C and S8C). Beyond
360 sampling the rare secondary structural changes, seen in the first
361 four examples, here MELD visits states separated by variations
362 in tertiary structure at the protein-protein interfaces (Fig. S9).
363 Therefore, starting with an ensemble of structures generated to
364 resolve 6RET, the inter-state hopping promoted by MELD's
365 enhanced sampling of the interface contacts (53), and refinement
366 by ReMDFF allowed for the resolution of three more conformations
367 of ATP synthase consistent with 6RDQ, 6RDK and 6RDW
368 (Tables: S6 and S7).

369 A key biophysical outcome that we make from the CryoFold
370 ensembles of ATP synthase is the flexibility of this motor's
371 peripheral stalk domains. Specifically, the OSCP hinge (chain
372

P) assumes a number of distinct open and closed conformations
373 with an RMSD of 3.3-6.4 Å (Fig. 5D) relative to the hinge from
374 6RET. The elastic coupling in ATP synthase has remained a topic
375 of contention in the bioenergy community with crystallographers
376 claiming minimum flexibility of the stalk regions (54), in sharp
377 contrast to single-molecule observations of "power-strokes" that
378 originate from deformations of the stalk (55). Within the
379 CryoFold ensembles incorporating all the states I-IV, we see that
380 the central stalk is in fact less flexible than the peripheral stalk
381 with an RMSD ranging between 2.4-3.8 Å relative to 6RET. So,
382 our results show that most of the elastic coupling in *polytomella*
383 ATP synthase comes from the peripheral stalk, rather than the
384 central stalk.
385

386 **F. Tests on soluble and membrane domains of a large ion channel with mid-resolution data.**
387 A second major challenge in *de novo* structure determination
388 arises from the modeling of complete transmembrane protein
389 systems, including structure of both the soluble and TM domains.
390 The refinement becomes particularly daunting for CryoFold, as
391 MELD simulations fail to capture structural changes from explicit
392 protein-membrane interactions (44). Consequently, the accuracy
393 of the model will depend on the structural information available
394 from the map, and less on the fidelity of the physical interactions
395 that underscore MELD.
396

397 Addressing this challenge, CryoFold was employed to model a
398 monomer from the pentameric Magnesium channel CorA, containing
399 349 residues, at 3.80 Å resolution(56) (pdb id: 3JCF,
400 EMDatabank: EMD-6551) (Figs. 6) and S10. An initial
401 topological prediction of the channel was obtained by flexibly
402 fitting of a linear polypeptide onto the $C\alpha$ trace obtained
403 from the cryo-EM density using MAINMAST. These traces were
404 already within 6.0 Å of the target $C\alpha$ conformation in 3JCF,
405 providing high-confidence coarse-grained information for MELD
406 to operate. Leveraging the MAINMAST trace, MELD was used to
407 perform local conformational sampling, regenerating most of the
408 secondary structures. The model with the highest cross-correlation
409 to the map was then refined using ReMDFF, finally resulting in
410 models which were at 2.90 Å RMSD to the native state. Even though
411 this model possessed high secondary structure content of 76%,
412 substantial unstructured regions remained both in the cytoplasmic
413 and the transmembrane regions, warranting a further round of
414 refinement. In the subsequent MELD-ReMDFF iteration, the
415 resulting models were 2.60 Å to the native state and agreed well
416 with the map with a CC of 0.84. Moreover, the CryoFold models
417 were comparable in geometry to that deposited in the database
418 (Fig. 6).
419

420 Discussion

421 The systems presented here have been chosen as challenging
422 problems to the methods that constitute CryoFold. We have not
423 over-optimized any aspect of the protocol to fit one problem,
424 rather complemented the uncertainties and weakness of one
425 method with the strengths of another. This approach is akin to
426 the consensus methods that are known to improve performance
427 over single methods in blind prediction challenges (57). A
428 selected combination of methods within CryoFold's plug-and-play
429 protocol will enable the prediction of completely unseen data
430 sets (Fig. S11), where the individual methods will potentially
431 fail.

432 While CryoFold appears promising for obtaining biomolecu- 493
 433 lar structures from cryo-EM, we are aware of some limitations. 494
 434 First, its success depends upon the correctness of the initial 495
 435 trace generated by MAINMAST. It is not clear when and 496
 436 whether the MD tools can recover from a wrong chain trace, 497
 437 particularly for resolving the transmembrane systems. Unlike 498
 438 Flpp3, repeating the CorA refinement with a poor-quality 499
 439 MAINMAST trace resulted in unreliable models. We do not
 440 have a good implicit membrane model to use in the MELD
 441 simulations and the use of explicit solvent would require many
 442 replicas, seeking more resources than currently available. Thus,
 443 by relying solely on the information coming from the density
 444 map we impose positional restraints and focus sampling on
 445 the transmembrane domains. Second, as with any MD sim-
 446 ulation of biomolecules, the force fields are still not perfect
 447 and larger structures will be a challenge for the searching and
 448 sampling, even with an accelerator such as MELD. Finally,
 449 in our current approach, MELD is the most computationally
 450 limiting, requiring between one and ten days of sampling with
 451 30 GPUs for the systems studied. This computational expense
 452 is not prohibitive using supercomputing resources available to
 453 academic researchers.

454 Despite the aforementioned limitations, CryoFold has been
 455 compared to the popular Rosetta protocols for TRPV1, ATP
 456 synthase and CorA. While for TRPV1 and CorA, Rosetta
 457 converged to models with unphysical overlap between the
 458 β -sheets (Fig. S5 and S12), a multi-protein refinement for
 459 ATP synthase could not be reproduced in ROSETTA-ES using
 460 standard resources, though individual chain refinements were
 461 achieved and are reported in Fig. S13. Thus, barring the Flpp3
 462 case at 1.8 Å, CryoFold was always found to offer higher quality
 463 models, but more importantly a diverse range of structures
 464 consistent with the expected biophysics.

465 A key benefit of this work is the ability to capture ensem- 521
 466 bles rather than single structures. Consequently, we identify 522
 467 conformations that are close to the native structure but also 523
 468 some alternative meta-stable states that are favored by the 524
 469 combination of force field and data. An important question 525
 470 follows – are these structures really relevant or just spurious? 526
 471 To this end, we have now validated using NMR and cryo-EM 527
 472 experiments that in addition to the narrow set of models con- 528
 473 sistent with one electron density map, there exists orthogonal 529
 474 states that are observed both in the experiments in CryoFold 530
 475 refinements. These orthogonal structures sampled by MELD 531
 476 are indeed leveraged in biological functions, as we shown by 532
 477 the open→close transition in Flpp3 or flexibility of the periph- 533
 478 eral stalks in elastic coupling of the ATP synthase example, 534
 479 yet behooves resolution by the limited sampling capacity of 535
 480 brute-force MD or MC sampling used in stationary structure 536
 481 determination. 537

482 Finally, evident from the 2016 and 2019 EMDB competition 538
 483 results, heterogeneous map resolutions affect the completeness 539
 484 of all the ensuing models. While a significant number of 540
 485 modelers prefer to truncate the more dynamic regions, MDFF 541
 486 offers a way to quantify uncertainty of the dynamic regions 542
 487 with root mean square deviations from an average model 543
 488 (27), and to correlate the inherent flexibility of proteins with 544
 489 the local resolution of density maps. Now, inside CryoFold,
 490 the fluid-like regions are even more thoroughly sampled by
 491 MELD offering the possibility of seeking hidden states in these
 492 fuzzy regions. Altogether, we present the first MD based

methodology for data-guided protein folding and ensemble
 refinement, bridging the strengths from two distinct areas
 of Biophysics. The implementation is semi-automated, and
 manual fitting is completely avoided. However, the user will
 require to control the I/O between the three methods, and
 optimize the default parameters as required. We have provided
 a GUI to facilitate this stage.

500 Conclusions

501 Structures, dynamics and function are interlinked. We often
 502 concentrate on a set of tools to determine structures from data
 503 and then use alternate computational techniques to determine
 504 dynamics between these metastable structures to ultimately
 505 elucidate biological functions. By leveraging the parallel algo-
 506 rithms with techniques such as CryoEM that capture multiple
 507 states (but an unknown number of them) tools that can go be-
 508 yond single structures to establish molecular dynamics directly
 509 from data. CryoFold is a first step in that direction.

510 Methods

511 The data-guided fold and fitting paradigm presented herein
 512 combines three real-space refinement methodologies, namely
 513 MELD, MAINMAST and ReMDFF. In what follows, these
 514 three formulations are articulated individually and the readers
 515 are referred to the original publications for details. Then,
 516 we outline the hybridization of the methods to provide a
 517 molecular dynamics-based *de novo* structure determination
 518 tool, CryoFold. Details of the setup for each individual system
 519 is outlined in Supplementary Information to showcase the
 520 different contexts in which CryoFold can operate.

521 **MELD:** Modeling Employing Limited Data (MELD) employs
 522 a Bayesian inference approach (eq. Eq. (1)) to incorporate em-
 523 pirical data into MD simulations(15, 28). The bayesian prior
 524 $p(\vec{x})$ comes from an atomistic force field (ff14SB sidechain,
 525 ff99SB backbone) and an implicit solvent model (Generalized
 526 born with neck correction, gb-neck2) (37, 38). The likelihood
 527 $p(\vec{D}|\vec{x})$, representing a bias towards known information, de-
 528 termines how well do the sampled conformations agree with
 529 known data, D . $p(\vec{D})$ refers to the likelihood of the data,
 530 which we take as a normalization term that can typically be
 531 ignored. Taken together,

$$532 \overbrace{p(\vec{x}|\vec{D})}^{\text{posterior}} = \frac{p(\vec{D}|\vec{x})p(\vec{x})}{p(\vec{D})} \sim \overbrace{p(\vec{D}|\vec{x})}^{\text{likelihood}} \overbrace{p(\vec{x})}^{\text{prior}}. \quad [1] \quad 533$$

534 MELD is designed to handle data with one or more of these
 535 features: sparsity, noise and ambiguity. Brute-force use of
 536 such data leads to incorrect models(58) as not all the data
 537 is compatible with the native state. MELD addresses the
 538 refinement of low-resolution data by enforcing only a fraction
 539 ($x\%$) of this data at every step of the MD simulation. Although
 540 x is kept fixed, the subset of data chosen to bias the simulation
 541 keeps changing with the simulation steps in a deterministic
 542 way. For a give nstructure all the data is evaluated, sorted
 543 according to their energy penalty and the $x\%$ with lowest
 544 energy guide the simulation until the next step. The data
 is incorporated as flat-bottom harmonic restraints $E(r_{ij})$ for

545 evaluating the likelihood ($p(\vec{D}|\vec{x})$).

$$546 \quad E(r_{ij}) = \begin{cases} \frac{1}{2}k(r_1 - r_2)(2r_{ij} - r_1 - r_2) & \text{if } r_{ij} < r_1 \\ \frac{1}{2}k(r_{ij} - r_2)^2 & \text{if } r_1 \leq r_{ij} < r_2 \\ 0 & \text{if } r_2 \leq r_{ij} < r_3 \\ \frac{1}{2}k(r_{ij} - r_3)^2 & \text{if } r_3 \leq r_{ij} < r_4 \\ \frac{1}{2}k(r_4 - r_3)(2r_{ij} - r_4 - r_3) & \text{if } r_4 \leq r_{ij}, \end{cases} \quad [2]$$

547 When these restraints are satisfied they do not contribute to
548 the energy or forces, contributing for flat bottom region of eq.
549 2 and (Fig. S12). When the restraints are not satisfied they
550 add energy penalties and force biases to the system – guiding it
551 to regions that satisfy a subset of the data, or conformational
552 envelopes. Details of MELD implementation are provided in
553 **Supplementary methods: Description of MELD.**

554 **MAINMAST.** MAINchain Model trAcing from Spanning Tree
555 (MAINMAST) is a *de novo* modeling program that directly
556 builds protein main-chain structures from an EM map of
557 around 4-5 Å or better resolutions(26). MAINMAST auto-
558 matically recognized main-chain positions in a map as dense
559 regions and does not use any known structures or structural
560 fragments. The procedure of MAINMAST consists of mainly
561 four steps (Fig. S14). In the first step, MAINMAST identifies
562 local dense points (LDPs) in an EM map by mean shifting
563 algorithm. All grid points in the map are iteratively shifted
564 by a gaussian kernel function and then merged to the clusters.
565 The representative points in the clusters are called LDPs. In
566 the second step, all the LDPs are connected by constructing a
567 minimum spanning tree (MST). It is found that the most edges
568 in the MST covers the main-chain of the protein structure
569 in EM map(26). In the third step, the initial tree structure
570 (MST) is refined iteratively by the so-called tabu search algo-
571 rithm. This algorithm attempts to explore a large search space
572 by using a list of moves that are recently considered and then
573 forbidden. In the final step, the longest path of the refined
574 tree is aligned with the amino acid sequence of the target pro-
575 tein. This process assigns optimal C α positions of the target
576 protein on the path and evaluates the fit of the amino acid
577 sequence to the longest path in a tree. Details of MAINMAST
578 implementation are provided in **Supplementary methods:**
579 **Description of MAINMAST.**

580 **Traditional MDFF.** The protocol for molecular dynamics flex-
581 ible fitting (MDFF) has been described in detail(6). Briefly,
582 a potential map V_{EM} is generated from the cryo-EM density
583 map, given by

$$584 \quad V_{EM}(\mathbf{r}) = \begin{cases} \zeta \left(1 - \frac{\Phi(\mathbf{r}) - \Phi_{thr}}{\Phi_{max} - \Phi_{thr}}\right) & \text{if } \Phi(\mathbf{r}) \geq \Phi_{thr}, \\ \zeta & \text{if } \Phi(\mathbf{r}) < \Phi_{thr}. \end{cases} \quad [3]$$

585 where $\Phi(\mathbf{r})$ is the biasing potential of the EM map at a point
586 \mathbf{r} , ζ is a scaling factor that controls the strength of the cou-
587 pling of atoms to the MDFF potential, Φ_{thr} is a threshold for
588 disregarding noise, and $\Phi_{max} = \max(\Phi(\mathbf{r}))$.

589 A search model is refined employing MD, where the tradi-
590 tional potential energy surface is modified by V_{EM} . The
591 density-weighted MD potential conforms the model to the
592 EM map, while simultaneously following constraints from
593 the traditional force fields. The output structure offers a
594 real-space solution, resolving the density with atomistically
595 detailed structures.

ReMDFF. While traditional MDFF works well with low- 596
597 resolution density maps, recent high-resolution EM maps have
598 proven to be more challenging. This is because high-resolution
599 maps run the risk of trapping the search model in a local
600 minimum of the density features. To overcome this unphysical
601 entrapment, resolution exchange MDFF (ReMDFF) employs
602 a series of MD simulations. Starting with $i = 1$, the i th map
603 in the series is obtained by applying a Gaussian blur of width
604 σ_i to the original density map. Each successive map in the
605 sequence $i = 1, 2, \dots, L$ has a lower σ_i (higher resolution),
606 where L is the total number of maps in the series ($\sigma_L = 0$ Å).
607 The fitting protocol assumes a replica-exchange approach
608 described in details(27) and illustrated in Fig. S15. At regular
609 simulation intervals, replicas i and j , of coordinates \mathbf{x}_i and
610 \mathbf{x}_j and fitting maps of blur widths σ_i and σ_j , are compared
611 energetically and exchanged with Metropolis acceptance
612 probability

$$613 \quad p(\mathbf{x}_i, \sigma_i, \mathbf{x}_j, \sigma_j) = \min \left(1, \exp \left(\frac{-U(\mathbf{x}_i, \sigma_j) - U(\mathbf{x}_j, \sigma_i) + U(\mathbf{x}_i, \sigma_i) + U(\mathbf{x}_j, \sigma_j)}{k_B T} \right) \right) \quad [4]$$

614 where k_B is the Boltzmann constant, $U(\mathbf{x}, \sigma)$ is the instan-
615 taneous total energy of the configuration \mathbf{x} within a fitting
616 potential map of blur width σ . Thus, ReMDFF fits the search
617 model to an initially large and ergodic conformational space
618 that is shrinking over the course of the simulation towards the
619 highly corrugated space described by the original MDFF poten-
620 tial map. Details of ReMDFF implementation are provided in
621 **Supplementary methods: Description of Resolution**
622 **exchange MDFF.**

CryoFold (MELD-MAINMAST-ReMDFF) protocol. Illustrated in 623
624 Fig. 1, the CryoFold protocol begins with MELD compu-
625 tations, which guided by backbone traces from MAINMAST
626 yields folded models. These models are flexibly fitted into
627 the EM density by ReMDFF to generate refined atomistic
628 structures.

- 629 1. First, information for the construction of Bayesian like-
630 lihood is derived from secondary structure predictions
631 (PSIPRED), which were enforced with a 70% confidence.
632 This percentage of confidence offers an optimal condition
633 for MELD to recover from the uncertainties in secondary
634 structure predictions(29). For membrane proteins, this
635 number can be increased to 80% when the transmembrane
636 motifs are well-defined helices. MELD extracts additional
637 prior information from the MD force field and the implicit
638 solvent model (see eq.1).
- 639 2. In the second step, any region determined with high
640 accuracy will be kept in place with cartesian restraints
641 imposed on the C α during the MELD simulations. This
642 way, the already resolved residues can fluctuate about
643 their initial position.
- 644 3. In the third step, distance restraints (e.g. from the C α
645 traces of MAINMAST) are derived. The application of
646 MAINMAST allows construction of pairwise interactions
647 as MELD-restraints directly from the EM density fea-
648 tures. Together with the cartesian restraints of step 2,
649 these MAINMAST-guided distance restraints are enforced
650 via flat-bottom harmonic potentials (see eq. 2) to guide
651 the sampling of a search model; notably, the search model

652 is either a random coil or manifests some topological fea-
653 tures when created by fitting the coil to the C α trace
654 with targeted MD. Depending upon the stage of CryoFold
655 refinements, only a percent of the cartesian and distance
656 restraints need be satisfied. The cartesian restraints are
657 often localized on the structured regions, while the dis-
658 tance restraints typically involve regions that are more
659 uncertain (e.g loop residues).

- 660 4. Fourth, a Temperature and Hamiltonian replica exchange
661 protocol (H,T-REMD) is employed to accelerate the sam-
662 pling of low-energy conformations in MELD(15, 28), re-
663 fining the secondary-structure content of the model. The
664 Hamiltonian is changed by changing the force constant
665 applied to the restraints. Simulations at higher replica
666 indexes have higher temperatures and lower (vanishing)
667 force constants so sampling is improved. At low replica
668 index, temperatures are low and the force constants are
669 enforced at their maximum value (but only a certain per-
670 cent of the restraints, the ones with lower energy, are
671 enforced). See SI for details for individual applications.
- 672 5. Fifth, cross-correlation of the H,T-REMD-generated struc-
673 tures with the EM-density is employed as a metric to select
674 the best model for subsequent refinement by ReMDFF
675 (Fig. S16). Resolution exchange across 5 to 11 maps with
676 successively increasing Gaussian blur of 0.5 Å (σ in eq.
677 4) sufficed to improve the cross-correlation and structural
678 statistics. The model with the highest EMringer score
679 forms the starting point of the next round of MELD sim-
680 ulations. Thereafter, another round ReMDFF is initiated,
681 and this iterative MELD-ReMDFF protocol continues
682 until the δ CC between two consecutive iterations is <0.1.

683 Throughout different rounds of iterative refinement, the struc-
684 tures from ReMDFF are used as seeds in new MELD simula-
685 tions. At the same time, distance restraints from the ReMDFF
686 model are updated and the pairs of residues present in those
687 interactions are enforced at different accuracy levels. As ex-
688 pected, the more rounds of refinement we do, the higher the
689 accuracy levels for the contacts is achieved in CryoFold. In
690 going through this procedure, the ensembles produced get
691 progressively narrower as we increase the amount of restraints
692 enforced. A video tutorial and the description of this implemen-
693 tation is provided in **Supplementary methods: Graphical**
694 **User Interface.**

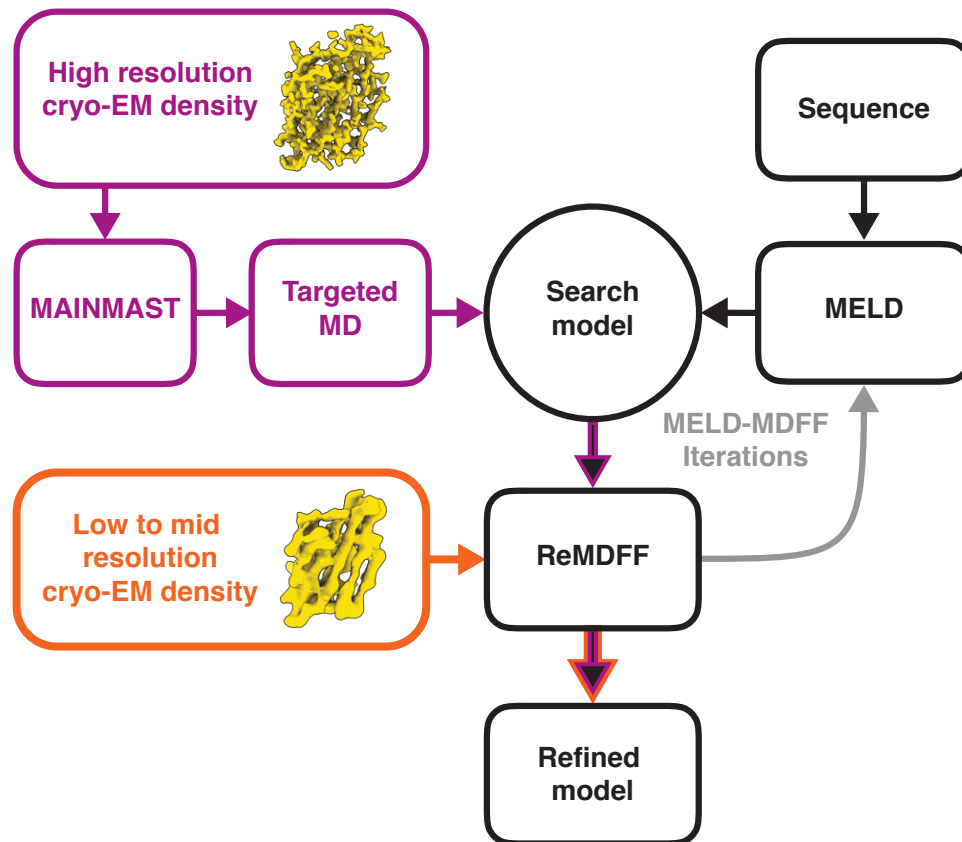


Fig. 1. An overview of the CryoFold protocol. For a high-resolution density map (data-rich case), backbone tracing is performed using MAINMAST to determine C α positions, and a random coil is fitted to these positions using targeted MD. This fitted protein model is subjected to the next MELD-ReMDFF cycles as a search model. For a low or medium resolution density map (data-poor case), a search model is constructed from primary sequence using MELD. This search model is fitted into the electron density using ReMDFF. The ReMDFF output is fed back to MELD for the next iteration, and the cycle continues until convergence.

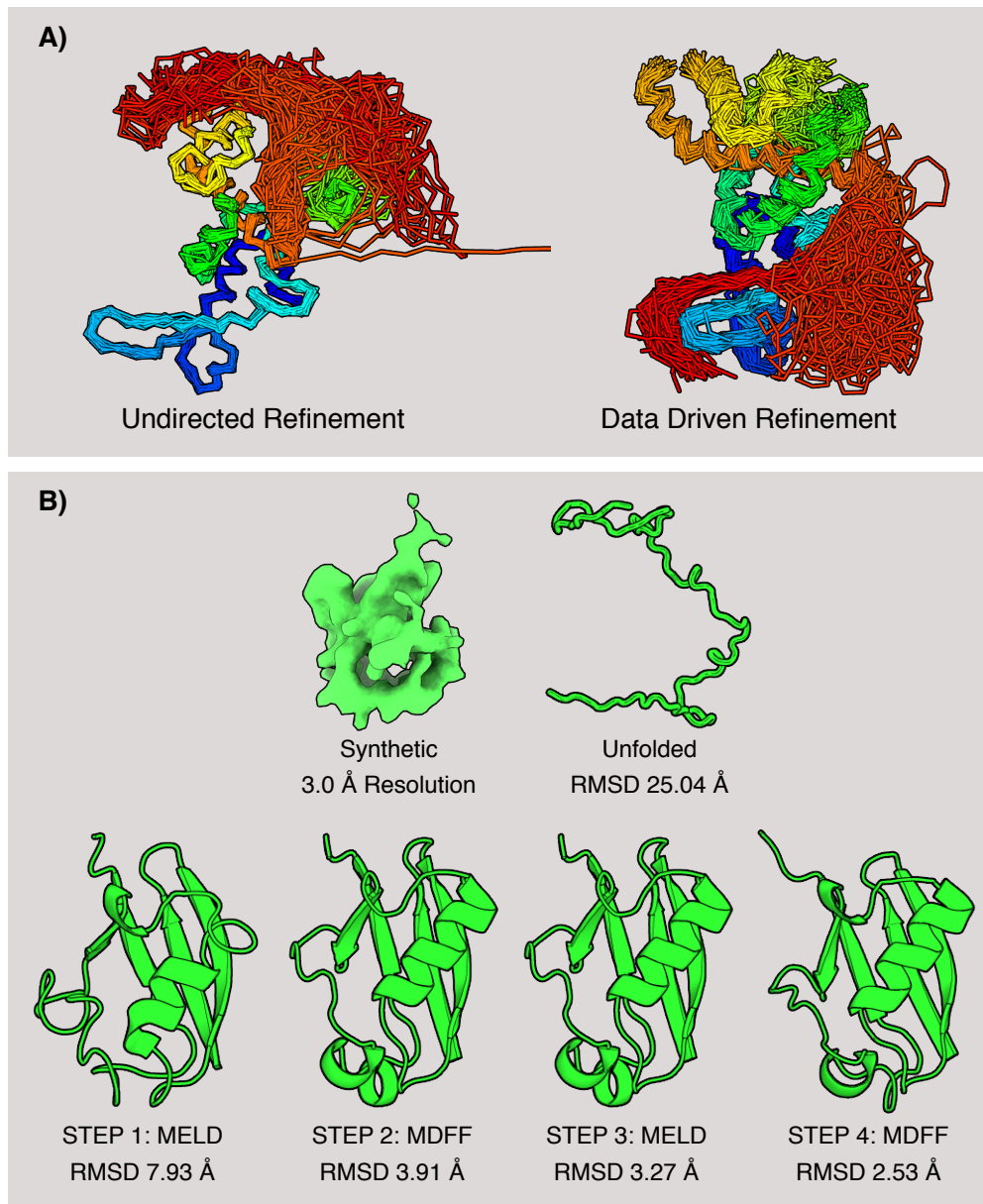


Fig. 2. Ensemble models for TRPV1 and the refinement protocol for ubiquitin. (A) Ensemble refinement with CryoFold showcased for the soluble domain of TRPV1. Several conformations from the TRPV1 ensemble are superimposed; color coding from blue (N-terminal) to red (C-terminal). In a MELD-only simulation, a soluble loop (indicated in red) artifactually interacted with the transmembrane domains. Following the data-guidance from ReMDFF, this loop interacted with the soluble domains and a more focused ensemble is derived that agrees with the electron density. **(B)** Stages of the refinement protocol for a test case, ubiquitin. The initial model is an unfolded coil. MELD was used to generate 50 search models from just the amino acid sequence, and no usage of the electron density data. Then, these models were rigid-fitted into the electron density using Chimera(59), and ranked based on their global cross-correlation. ReMDFF refined the best rigid-fitted model even further. The ReMDFF model with the highest Cross Correlation (CC) to the density map served as a template for the subsequent iteration with MELD. In two consecutive MELD-ReMDFF iterations the RMSD of the folded model relative to the crystal structure (1UBQ) attenuated from 25.04 Å to 2.53 Å

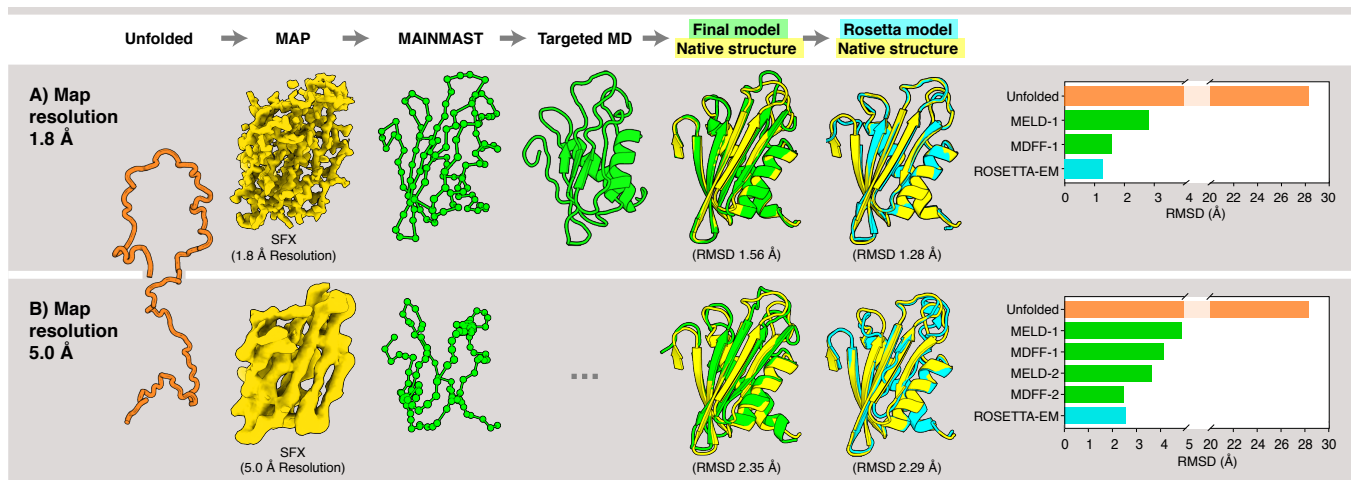


Fig. 3. Hybrid structure determination of Flpp3. (A) High-resolution density map at 1.8 Å resolution. An unfolded structure was used as the initial model. A SFX density map at 1.8 Å resolution was employed to generate the C α position (green spheres) using MAINMAST, and the initial model was fitted into these positions by targeted MD. The resulting structure (green cartoon model) was then subjected to MELD-ReMDFF refinement. This procedure yielded a structure with RMSD of 1.56 Å relative to the native SFX structure (yellow). The Rosetta-EM model (cyan) has an RMSD of 1.28 Å with respect to the SFX structure. (B) Lower-resolution density map at 5 Å resolution. An initial C α trace in the map was computed using MAINMAST. Subsequent MELD-ReMDFF refinement resulted in a structure (green cartoon model) with an RMSD of 2.29 Å from the SFX structure (yellow). The best Rosetta-EM model has (cyan) an RMSD of 2.35 Å to the SFX structure. Barplots depict the evolution of RMSD of the CryoFold models with each subsequent MELD-ReMDFF refinement.

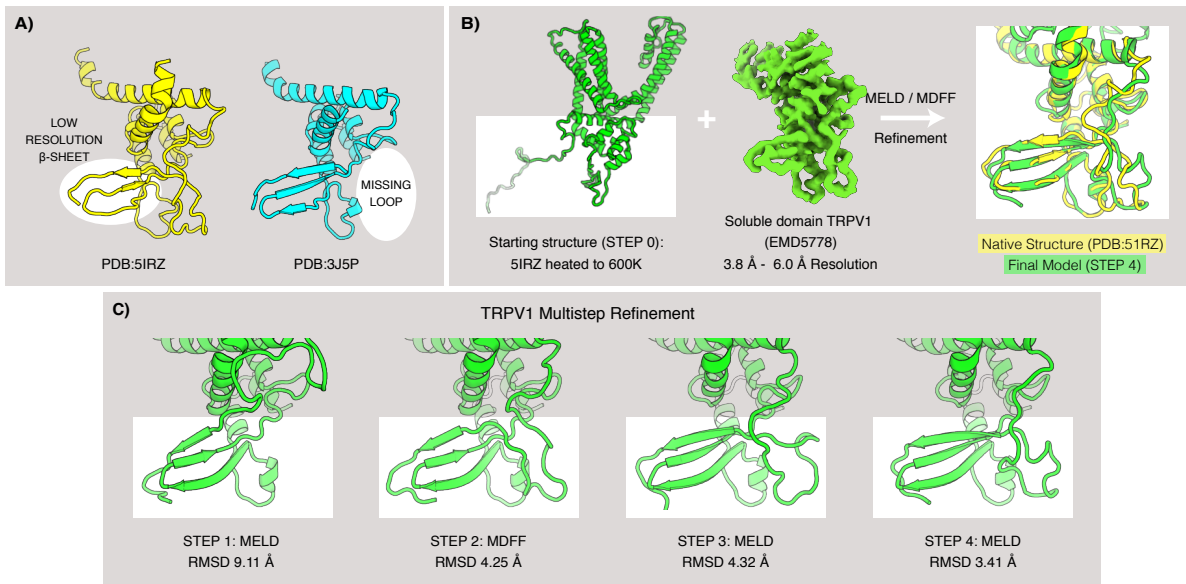


Fig. 4. Modeling of the soluble domain of TRPV1. (A) TRPV1 structures deposited in 2016 (pdb 5IRZ in yellow) and in 2013 (pdb 3J5P in cyan in cartoon representation, showing the latter has a more resolved β -sheet while the former possess an additional extended loop. (B) The 5IRZ model was heated at 600 K using brute-force MD, while constraining the α helices. After 10 ns of simulation, this treatment resulted in a search model with the loop regions significantly deviated and the β sheets completely denatured. The search model was subjected to MELD-ReMDFF refinement. A single round of MELD regenerated most of the β -sheet from this random chain, however the 5- to 15-residue long interconnecting loops still occupied non-native positions. Subsequent ReMDFF refinement with the 5IRZ density resurrected the loop positions. One more round of the MELD and ReMDFF resulted in the further refinement of the model. The final refined model agrees well with 5IRZ (C) Progress of the refinement in each step of CryoFold. MELD step 1 shows the β sheets modeled correctly, while the loops recovered in MDFF step 2, and refinement was complete by step 4.

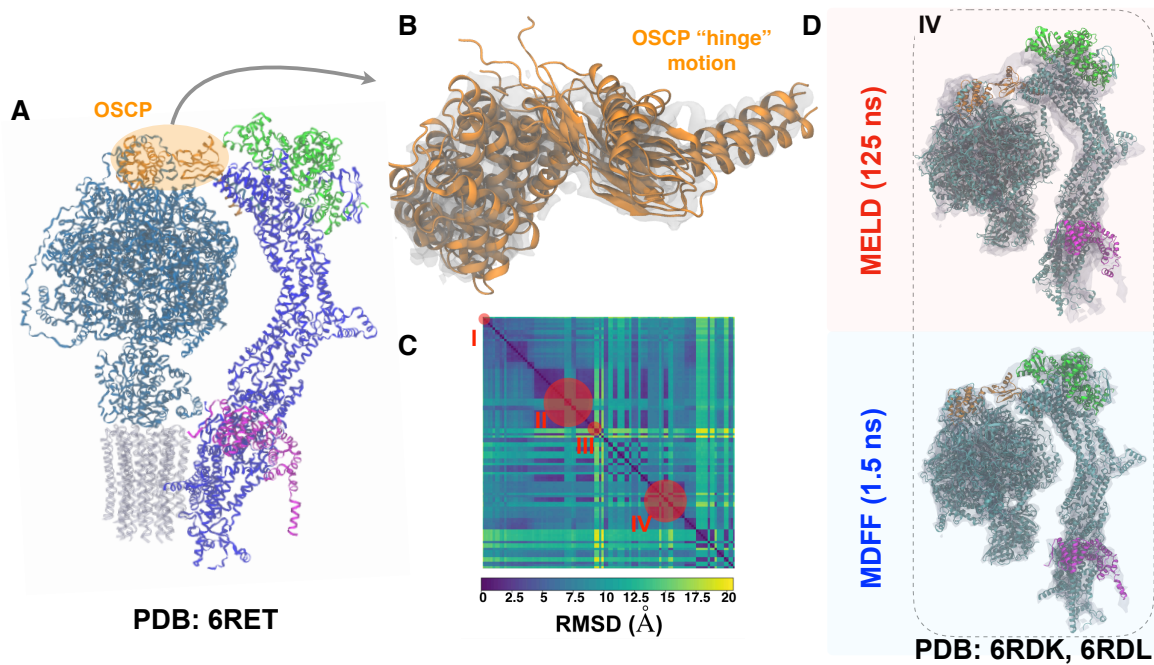


Fig. 5. **CryoFold samples several biologically relevant states of the soluble domain of mitochondrial F₁ - F₀ ATP-synthase.** We modeled mitochondrial F₁ - F₀ ATPsynthase starting from pdb 6RET (state I) and excluding the grey region embedded in the membrane from refinement. CryoFold samples different conformations through a hinge motion in the OSCP region (orange) connecting the arm (blue) with the rotary domains (cyan). Clustering and 2D-RMSD analysis shows Cryofold samples conformations of additional ATPsynthase states represented by pdb codes 6RDK, 6RDL (state IV). Other states represented by pdb codes 6RDQ, 6RDR (state II) and 6RDW, 6RDX (state III) are included in SI.

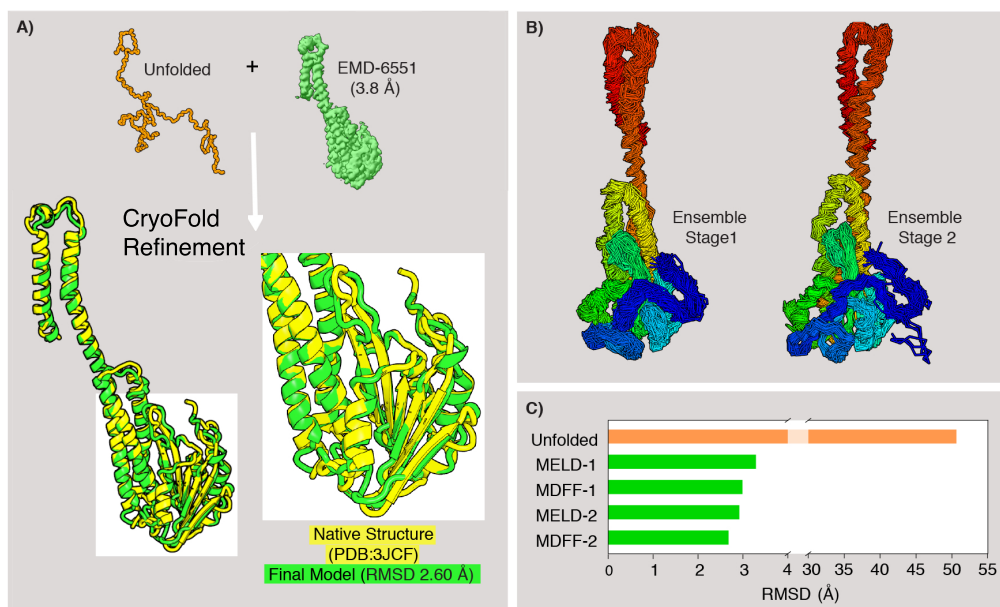


Fig. 6. Modeling transmembrane Magnesium-channel CorA. (A) The CryoFold protocol on CorA. A starts from an $C\alpha$ trace based Cryo-EM density map using MAINMAST and refined through different cycles of MELD and MDFF produces a structure that agrees extremely well with the native structure (yellow), featuring accurate beta structures. (B) CryoFold produces narrower, more constraint ensembles as we iterate through MELD/MDFF. (C) The evolution of the RMSD of CryoFold models with each MELD-ReMDFF refinement. The end-model is 2.60 Å RMSD from the native structure.

695 **ACKNOWLEDGMENTS.** AS and CG acknowledge start-up funds
696 from the SMS and CASD at Arizona State University, CAREER
697 award by NSF-MCB 1942763 and the resources of the OLCF at
698 the Oak Ridge National Laboratory, which is supported by the Of-
699 fice of Science at DOE under Contract No. DE-AC05-00OR22725,
700 made available via the INCITE program. ET laboratory is sup-
701 ported by NIH (P41GM104601); ET, AS and MS acknowledge
702 NIH (R01GM098243-02). This research is part of the Blue Wa-
703 ters sustained-petascale computing project, which is supported
704 by the National Science Foundation (Awards OCI-0725070 and
705 ACI-1238993) and the state of Illinois. KD and AP appreciate
706 support from a PRAC computer allocation supported by NSF
707 Award ACI1514873, support from NIH Grant GM125813 and the
708 Laufer Center. AP appreciates start-up support from the University
709 of Florida. DK acknowledges support from the National Insti-
710 tutes of Health (R01GM123055), the National Science Foundation
711 (DMS1614777, CMMI1825941), and the Purdue Institute of Drug
712 Discovery. WVH acknowledges NIH (R01GM112077).

713 References

714 1. P Afonine, J Headd, T Terwilliger, P Adams, New tool: phenix.real_space_refine. *Comput.*
715 *Crystallogr. NewsL* **4**, 43–44 (2013).
716 2. JC Burnett, JJ Rossi, RNA-based therapeutics: current progress and future prospects. *Chem.*
717 *Biol.* **19**, 60–71 (2012).
718 3. MP Rout, A Sali, Principles for Integrative Structural Biology Studies. *Cell* **177**, 1384–1403
719 (2019).
720 4. J Frank, A Ourmazd, Continuous changes in structure mapped by manifold embedding of
721 single-particle data in cryo-EM. *Methods* **100**, 61–67 (2016).
722 5. M Feig, Computational protein structure refinement: almost there, yet still so far to go. *Wiley*
723 *Interdiscip. Rev. Comput. Mol. Sci.* **7**, e1307 (2017).
724 6. LG Trabuco, E Villa, K Mitra, J Frank, K Schulten, Flexible fitting of atomic structures into
725 electron microscopy maps using molecular dynamics. *Structure* **16**, 673–683 (2008).
726 7. PG Wolynes, Folding funnels and energy landscapes of larger proteins within the capillarity
727 approximation. *Proc. Natl. Acad. Sci.* **94**, 6170–6175 (1997).
728 8. R Zhou, BJ Berne, R Germain, The free energy landscape for beta hairpin folding in explicit
729 water. *Proc Natl Acad Sci U S A* **98**, 14931–14936 (2001).
730 9. S Ovchinnikov, H Park, DE Kim, F DiMaio, D Baker, Protein structure prediction using rosetta
731 in casp12. *Proteins: Struct. Funct. Bioinforma.* **86**, 113–121 (2018).
732 10. CC Jolley, SA Wells, P Fromme, MF Thorpe, Fitting low-resolution cryo-EM maps of proteins
733 using constrained geometric simulations. *Biophys. J.* **94**, 1613–1621 (2008).
734 11. TA Hopf, et al., Three-dimensional structures of membrane proteins from genomic sequenc-
735 ing. *Cell* **149**, 1607–1621 (2012).
736 12. TA Hopf, et al., Sequence co-evolution gives 3d contacts and structures of protein complexes.
737 **3**, e03430 (2014).
738 13. A Sali, TL Blundell, Comparative protein modelling by satisfaction of spatial restraints. *J. Mol.*
739 *Biol.* **234**, 779 (1993).
740 14. N Eswar, et al., Comparative protein structure modeling using MODELLER. *Curr. Protoc.*
741 *Bioinf.*, 5–6 (2006).
742 15. A Perez, JL MacCallum, KA Dill, Accelerating molecular simulations of proteins using
743 Bayesian inference on weak information. *Proc. Natl. Acad. Sci. United States Am.* **112**,
744 11846–11851 (2015).
745 16. J Comer, et al., The Adaptive Biasing Force Method: Everything You Always Wanted To Know
746 but Were Afraid To Ask. *The J. Phys. Chem. B* **119**, 1129–1151 (2015).
747 17. M Bonomi, C Camilloni, A Cavalli, M Vendruscolo, MetaInference: A Bayesian inference
748 method for heterogeneous systems. *Sci. Adv.* **2**, e1501177 – e1501177 (2016).
749 18. DE Shaw, et al., Millisecond-scale molecular dynamics simulations on Anton in SC'09: *Pro-*
750 *ceedings of the Conference on High Performance Computing Networking, Storage and Anal-*
751 *ysis.* (ACM, New York, NY, USA), pp. 39:1–39:11 (2009).
752 19. JC Robertson, A Perez, K Dill, MELD MD Folds Nonthreadables, Giving Native Structures
753 and Populations. *J. chemical theory computation* **14**, 6734 – 6740 (2018).
754 20. W Ray, et al., De novo protein structure determination from near-atomic-resolution cryo-em
755 maps. *Nat. Methods* **12**, 335 (2015).
756 21. KT Simons, R Bonneau, I Ruczinski, D Baker, Ab initio protein structure prediction of CASP
757 III targets using ROSETTA. *Proteins: Struct. Funct. Bioinforma.* **37**, 171–176 (1999).
758 22. F DiMaio, A Leaver-Fay, P Bradley, D Baker, I Andr  , Modeling symmetric macromolecular
759 structures in rosetta3. *PLOS ONE* **6**, 1–13 (2011).
760 23. F DiMaio, et al., Atomic-accuracy models from 4.5   cryo-electron microscopy data with
761 density-guided iterative local refinement. *Nat. Methods* **12**, 361–365 (2015).
762 24. SP Leelananda, S Lindert, Iterative molecular dynamics–rosetta membrane protein structure
763 refinement guided by cryo-em densities. *J. Chem. Theory Comput.* **13**, 5131–5145 (2017).
764 25. S Wickles, et al., A structural model of the active ribosome-bound membrane protein inser-
765 tase YidC. *eLife* **3:e03035** (2014) (17 pages).
766 26. G Terashi, D Kihara, De novo main-chain modeling for EM maps using MAINMAST. *Nat.*
767 *Commun.* **9**, 1618 (2018).
768 27. A Singharoy, et al., Molecular dynamics-based refinement and validation for sub-5   cryo-
769 electron microscopy maps. *eLife* **10.7554/eLife.16105** (2016).
770 28. JL MacCallum, A Perez, KA Dill, Determining protein structures by combining semireliable
771 data with atomistic physical models by Bayesian inference. *Proc. Natl. Acad. Sci. United*
772 *States Am.* **112**, 6985–6990 (2015).

29. SR Jones, et al., Loss of autoreceptor functions in mice lacking the dopamine transporter. *773*
Nat. Neurosci. **2**, 649–655 (1999). *774*
30. AM Karmali, TL Blundell, N Furnham, Model-building strategies for low-resolution X-ray crys- *775*
tallographic data. *Acta Cryst. D* **65**, 121–127 (2009). *776*
31. R Baradaran, JM Berrisford, GS Minhas, LA Sazanov, Crystal structure of the entire respira- *777*
tory complex I. *Nature* **494**, 443–448 (2013). *778*
32. S Piana, K Lindorff-Larsen, DE Shaw, Atomic-level description of ubiquitin folding. *Proc. Natl.* *779*
Acad. Sci. **110**, 5915 – 5920 (2013). *780*
33. S Vijay-Kumar, C Bugg, W Cook, Structure of ubiquitin refined at 1.8   resolution. *J. Mol. Biol.* *781*
194, 531–544 (1987). *782*
34. J Zook, et al., NMR Structure of Francisella tularensis Virulence Determinant Reveals Structural *783*
Homology to Bet v1 Allergen Proteins. *Struct. (London, Engl. : 1993)* **23**, 1116–1122 *784*
(2015). *785*
35. J Zook, et al., XFEL and NMR Structures of Francisella Lipoprotein Reveal Conformational *786*
Space of Drug Target against Tularemia. *Structure* **28**, 540–547.e3 (2020). *787*
36. SP Leelananda, S Lindert, Using nmr chemical shifts and cryo-em density restraints in iter- *788*
ative rosetta-md protein structure refinement. *J. Chem. Inf. Model.* **60**, 2522–2532 (2020) *789*
PMID: 31872764. *790*
37. JA Maier, et al., ff14SB: Improving the Accuracy of Protein Side Chain and Backbone Param- *791*
eters from ff99SB. *J. Chem. Theory Comput.* **11**, 3696–3713 (2015). *792*
38. H Nguyen, DR Roe, C Simmerling, Improved Generalized Born Solvent Model Parameters *793*
for Protein Simulations. *J. Chem. Theory Comput.* **9**, 2020–2034 (2013). *794*
39. J Huang, et al., CHARMM36m: an improved force field for folded and intrinsically disordered *795*
proteins. *Nat. Methods* **14**, 71–73 (2017). *796*
40. A Kucukelbir, FJ Sigworth, HD Tagare, Quantifying the local resolution of cryo-EM density *797*
maps. *Nat. Methods* **11**, 63–65 (2014). *798*
41. M Liao, E Cao, D Julius, Y Cheng, Structure of the TRPV1 ion channel determined by electron *799*
cryo-microscopy. *Nature* **504**, 107–112 (2013). *800*
42. A Kucukelbir, FJ Sigworth, HD Tagare, Quantifying the local resolution of cryo-em density *801*
maps. *Nat. methods* **11**, 63 (2014). *802*
43. Y Wang, et al., Constructing atomic structural models into cryo-em densities using molecular *803*
dynamics- pros and cons. *J. Struct. Biol.* **204**, 319 – 328 (2018). *804*
44. JL MacCallum, A Perez, KA Dill, Determining protein structures by combining semireliable *805*
data with atomistic physical models by Bayesian inference. *Proc. Natl. Acad. Sci.* **112**, 6985– *806*
6990 (2015). *807*
45. L Zubcevic, AL Hsu, MJ Borgnina, SY Lee, Symmetry transitions during gating of the trpv2 ion *808*
channel in lipid membranes. *eLife* **8**, e45779 (2019). *809*
46. JW Vant, et al., Flexible Fitting of Small Molecules into Electron Microscopy Maps Using *810*
Molecular Dynamics Simulations with Neural Network Potentials. *J. Chem. Inf. Model.* **60**, *811*
2591–2604 (2020). *812*
47. LA Abriata, MD Peraro, Will Cryo-Electron Microscopy Shift the Current Paradigm in Protein *813*
Structure Prediction? *J. chemical information modeling* **60**, 2443–2447 (2020). *814*
48. CL Lawson, et al., Outcomes of the 2019 emdataproject model challenge: validation of *815*
cryo-em models at near-atomic resolution. *bioRxiv* (2020). *816*
49. 2019 model metrics challenge | em validation challenges (<https://challenges.emdataproject.org/?q=model-metrics-challenge-2019>) (2019) (Accessed on 06/27/2020). *817*
50. CJ Williams, et al., MolProbity: More and better reference data for improved all-atom structure *818*
validation. *Protein Sci.* **27**, 293–315 (2018). *819*
51. F Guo, W Jiang, Single particle cryo-electron microscopy and 3-d reconstruction of viruses. *820*
Methods Mol Biol **1117**, 401–443 (2014). *821*
52. BJ Murphy, et al., Rotary substates of mitochondrial ATP synthase reveal the basis of flexible *822*
F1-Fo coupling. *Science* **364**, eaaw9128 (2019). *823*
53. JA Morrone, A Perez, J MacCallum, KA Dill, Computed binding of peptides to proteins with *824*
meld-accelerated molecular dynamics. *J. Chem. Theory Comput.* **13**, 870–876 (2017). *825*
54. JL Rubinstein, Structure of the mitochondrial ATP synthase by electron cryomicroscopy. *The* *826*
EMBO J. **22**, 6182–6192 (2003). *827*
55. JL Martin, R Ishmukhametov, D Spletzer, T Hornung, WD Frasch, Elastic coupling power *828*
stroke mechanism of the F1-ATPase molecular motor. *Proc. Natl. Acad. Sci.* **115**, 5750– *829*
5755 (2018). *830*
56. D Matthies, et al., Cryo-EM Structures of the Magnesium Channel CorA Reveal Symmetry *831*
Breaks upon Gating. *Cell* **164**, 747–756 (2016). *832*
57. E Wilson, G Hirneise, A Singharoy, KS Anderson, Total predicted mhc-i epitope load is in- *833*
versely associated with mortality from sars-cov-2. *medRxiv* (2020). *834*
58. BC Goh, et al., Computational methodologies for real-space structural refinement of large *835*
macromolecular complexes. *Annu. Rev. Biophys.* **45**, 253–278 (2016) PMID: 27145875. *836*
59. EF Pettersen, et al., UCSF Chimera - A visualization system for exploratory research and *837*
analysis. *J. Comp. Chem.* **25**, 1605–1612 (2004). *838*
839