



HAL
open science

SeaMoon: Prediction of molecular motions based on language models

Valentin Lombard, Dan Timsit, Sergei Grudinin, Elodie Laine

► **To cite this version:**

Valentin Lombard, Dan Timsit, Sergei Grudinin, Elodie Laine. SeaMoon: Prediction of molecular motions based on language models. 2024. hal-04801636

HAL Id: hal-04801636

<https://hal.science/hal-04801636v1>

Preprint submitted on 25 Nov 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 SeaMoon: Prediction of molecular motions based on
2 language models

3 Valentin Lombard¹, Dan Timsit¹, Sergei Grudinin^{*2}, Elodie Laine^{*1,3}

4 ¹ Sorbonne Université, CNRS, IBPS, Laboratoire de Biologie Computationnelle et
5 Quantitative (LCQB), 75005 Paris, France.

6 ² Université Grenoble Alpes, CNRS, Grenoble INP, LJK, 38000 Grenoble, France.

7 ³ Institut Universitaire de France (IUF).

8 * corresponding authors: sergei.grudinin@univ-grenoble-alpes.fr,
9 elodie.laine@sorbonne-universite.fr

Abstract

How protein move and deform determines their interactions with the environment and is thus of utmost importance for cellular functioning. Following the revolution in single protein 3D structure prediction, researchers have focused on repurposing or developing deep learning models for sampling alternative protein conformations. In this work, we explored whether continuous compact representations of protein motions could be predicted directly from protein sequences, without exploiting nor sampling protein structures. Our approach, called SeaMoon, leverages protein Language Model (pLM) embeddings as input to a lightweight ($\sim 1\text{M}$ trainable parameters) convolutional neural network. SeaMoon achieves a success rate of up to 40% when assessed against $\sim 1\,000$ collections of experimental conformations exhibiting a wide range of motions. SeaMoon capture motions not accessible to the normal mode analysis, an unsupervised physics-based method relying solely on a protein structure's 3D geometry, and generalises to proteins that do not have any detectable sequence similarity to the training set. SeaMoon is easily retrainable with novel or updated pLMs.

Keywords: protein motion, protein language models, transfer learning, PCA, deep learning

28 Introduction

29 Proteins coordinate and regulate all biological processes by adapting their 3D shapes
30 to their environment and cellular partners. Deciphering the complexities of how proteins
31 move and deform in solution is thus of utmost importance for understanding the cellular
32 machinery. Yet, despite spectacular advances in protein structure determination and pre-
33 diction, comprehending protein conformational heterogeneity remains challenging (Lane,
34 2023; Miller and Phillips, 2021; Henzler-Wildman and Kern, 2007).

35 Many recent approaches have concentrated on repurposing the protein structure predic-
36 tion neural network AlphaFold2 (Jumper et al., 2021) to generate conformational diversity
37 (Sala et al., 2023). Guiding the predictions with state-annotated templates proved suc-
38 cessful for modelling the multiple functional states of a couple of protein families (Faezov
39 and Dunbrack Jr, 2023; Heo and Feig, 2022). In addition, massive sampling strate-
40 gies have shown promising results for protein complexes (Wallner, 2023) (Wallner, 2023;
41 Johansson-Åkhe and Wallner, 2022) with notable success in the blind CASP15-CAPRI
42 assessment (Lensink et al., 2023). While they can be deployed seamlessly with parallelized
43 implementations (Brysbaert et al., 2024), they remain highly resource-intensive.

44 Other strategies have explored promoting diversity by modulating and disentangling
45 evolutionary signals (Sfriso et al., 2016). The rationale is that amino acid co-variations
46 in evolution reflect 3D structural constraints (Benner and Gerloff, 1991; Göbel et al.,
47 1994; Ortiz et al., 1999; Lapedes et al., 1999; Giraud et al., 1999; Thomas et al., 2005;
48 Weigt et al., 2009). These evolutionary patterns can be extracted directly from align-
49 ments of evolutionary related sequences, or, as shown more recently, by modeling raw
50 sequences at scale with protein language models (Bepler and Berger, 2021; Elnaggar
51 et al., 2022; Lin et al., 2023). Inputting shallow, masked, corrupted or sub-sampled
52 alignments to AlphaFold2 allowed for modelling distinct conformations for a few protein
53 families (Kalakoti and Wallner, 2024; Wayment-Steele et al., 2023; Del Alamo et al.,
54 2022; Stein and Mchaourab, 2022). Nevertheless, contradictory findings have highlighted
55 difficulties in rationalising the effectiveness of these modifications and interpreting them,
56 particularly for metamorphic proteins (Porter et al., 2024; Chakravarty and Porter, 2022;
57 Chakravarty et al., 2023).

58 More classically, physics-based molecular dynamics (MD) is a method of choice to
59 probe protein conformational landscapes (Hollingsworth and Dror, 2018). Nonetheless,
60 the time scales amenable to MD simulations on standard hardware remain much smaller
61 than those spanned by slow molecular processes (Chen et al., 2023). This limitation has
62 stimulated the development of hybrid approaches combining MD with machine learning
63 (ML) toward accelerating or enhancing sampling (Noé et al., 2020). Deep neural networks
64 can help to identify collective variables from MD simulations as part of importance-
65 sampling strategies (Chen et al., 2023; Belkacemi et al., 2021; Bonati et al., 2021; Wang
66 et al., 2020; Ribeiro et al., 2018). Or they may directly generate conformations according
67 to a probability distribution learnt from MD trajectories or sets of experimental structures
68 (Zheng et al., 2024; Lu et al., 2023; Ramaswamy et al., 2021; Noé et al., 2019). Diffusion-
69 based architectures (Abramson et al., 2024; Zheng et al., 2024; Jing et al., 2023) and
70 the more general flow-matching framework (Jing et al., 2024) provide highly efficient and
71 flexible means to generate diverse conformations conditioned on cellular partners and
72 ligands. Nevertheless, they are prone to hallucination, and models trained across protein

73 families still fail to approximate solution ensembles (Abramson et al., 2024).

74 On the other hand, the normal mode analysis (NMA) represents a data- and compute-
75 inexpensive unsupervised alternative for accessing large-scale, shape-changing protein mo-
76 tions (Hayward and Go, 1995). In particular, the NOLB method predicts protein func-
77 tional transitions in real-time by deforming single structures along a few collective coordi-
78 nates inferred with the NMA (Grudinin et al., 2020; Hoffmann and Grudinin, 2017). The
79 generated conformations are physically plausible and stereochemically realistic. However,
80 the results strongly depend on the 3D geometry of the starting structure, and although
81 some of the initial topological constraints can be easily alleviated (Laine and Grudinin,
82 2021), the NMA remains unsuitable for modelling extensive secondary structure rear-
83 rangements.

84 Training and benchmarking predictive methods is difficult due to the sparsity and
85 inhomogeneity of the available experimental data (Berman et al., 2000). X-ray crys-
86 tallography, cryogenic-electron microscopy (cryo-EM), and nuclear magnetic resonance
87 spectroscopy (NMR) have provided invaluable insights into protein diverse conformational
88 states (Ramelot et al., 2023; Miller and Phillips, 2021), but only for a relatively small num-
89 ber of proteins (Bryant, 2023). Small-angle X-ray or neutron scattering (SAXS, SANS)
90 and high-speed atomic force microscopy (HS-AFM) techniques allow for directly probing
91 continuous protein heterogeneity, but with limited structural resolution (Trehella, 2022;
92 Martel and Gabel, 2022; Flechsig and Ando, 2023).

93 Ongoing community-wide efforts aim at revealing the full potential of the available
94 structural data by collecting, clustering, curating, visualising and functionally annotating
95 experimental protein structures together with high-quality predicted models (Wankowicz
96 and Fraser, 2023; Ramelot et al., 2023; Ellaway et al., 2023; Varadi et al., 2022; Modi and
97 Dunbrack Jr, 2022; Parker et al., 2022; Tordai et al., 2022; Pándy-Szekeres et al., 2023).
98 For instance, the DANCE method produces movie-like visual narratives and compact con-
99 tinuous representations of protein conformational diversity, interpreted as *linear motions*,
100 from static 3D snapshots (Lombard et al., 2024). DANCE application to the Protein Data
101 Bank (PDB) (Berman et al., 2000) revealed that the conformations observed for most
102 protein families lie on a low-dimensional *manifold*. Classical dimensionality reduction
103 techniques can learn this manifold and generate unseen conformations with reasonable
104 accuracy, albeit only in close vicinity of the training set (Lombard et al., 2024).

105 Here, we explored the possibility of predicting protein motions directly from amino acid
106 sequences without exploiting nor sampling protein 3D structures. To do so, we lever-
107 aged protein Language Models (pLMs) pre-trained through self-supervision over large
108 databases of protein-related data. Our approach, SEAquencetoMOtioON or SeaMoon, is
109 a 1D convolutional neural network inputting a protein sequence pLM embedding and out-
110 putting a set of 3D displacement vectors (**Fig. 1**). The latter define protein residues' rela-
111 tive motion amplitudes and directions. We tested whether SeaMoon could capture the *lin-*
112 *ear motion manifold* underlying experimentally resolved conformations across thousands
113 of diverse protein families (Lombard et al., 2024). To this end, we devised an objective
114 function invariant to global translations, rotations, and dilatations in 3D space. SeaMoon
115 achieved a success rate similar to the normal mode analysis (NMA) when inputting purely
116 sequence-based pLM embeddings (Lin et al., 2023) without any knowledge about protein
117 3D structures. It could generalise to proteins without any detectable sequence similarity
118 to the training set and capture motions not directly accessible from protein 3D geometry.

119 Injecting implicit structural knowledge with sequence-structure bilingual or multimodal
120 pLMs (Hayes et al., 2024; Heinzinger et al., 2023) further boosted the performance. This
121 work establishes a community baseline and paves the way for developing evolutionary-
122 and physics-informed neural networks to predict continuous protein motions.

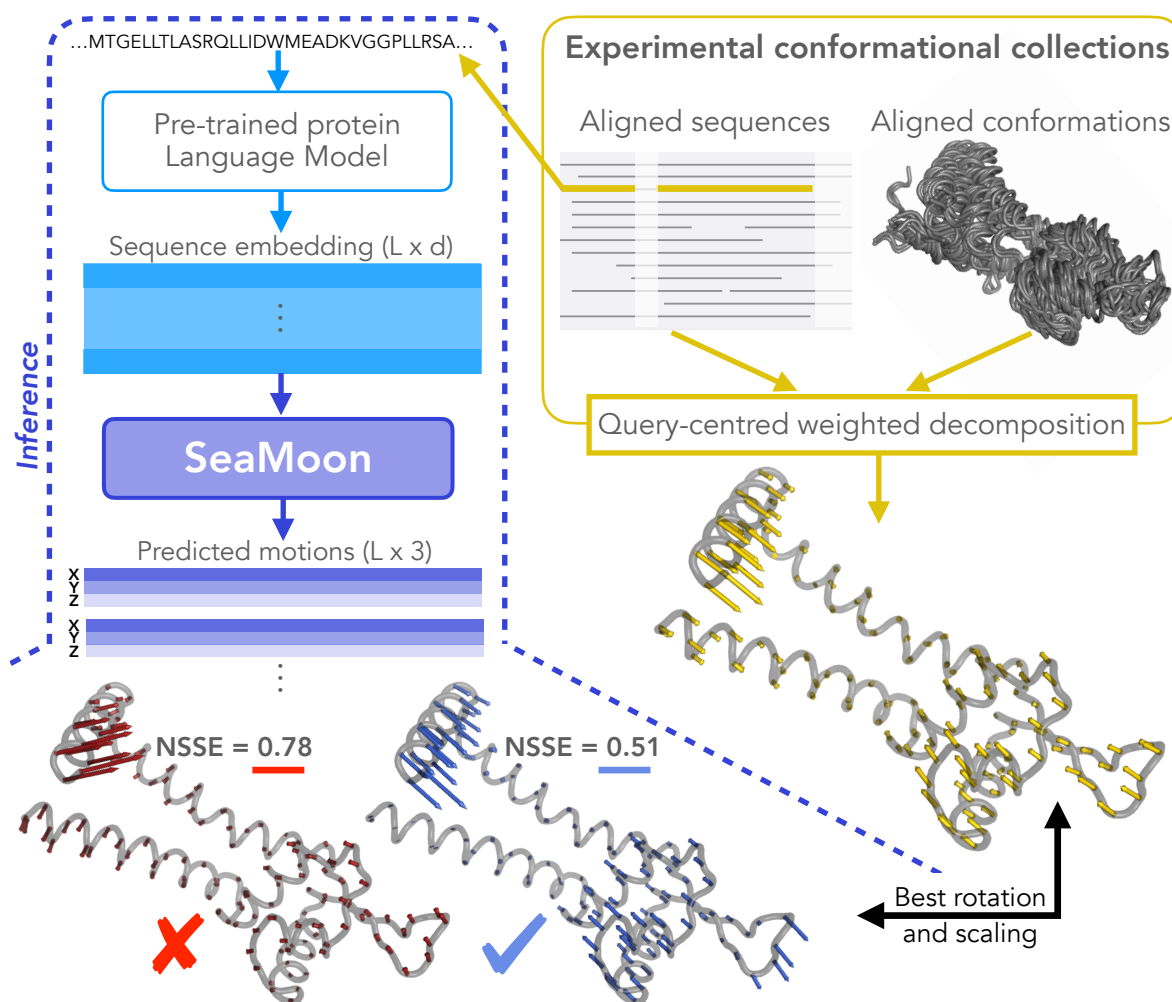


Figure 1: **Outline of SeaMoon’s approach.** SeaMoon takes as input a high-dimensional $L \times d$ matrix representation of a protein sequence of length L computed by a pre-trained pLM. It outputs a set of 3D vectors of length L representing linear motions. The training procedure regresses these output motions (blue and red arrows) against ground-truth ones (yellow arrows) extracted from experimental conformational collections through principal component analysis. For this, SeaMoon identifies the transformation (rotation and scaling) minimising their discrepancy, computed as a sum-of-squares error (SSE). We consider predictions with a normalised error (NSSE) smaller than 0.6 as acceptable. We show the query protein 3D structure only for illustrating the motions, it is not used by SeaMoon nor by the pLM generating the input embeddings..

123 Results and Discussion

124 The approach introduced in this work, SeaMoon, predicts continuous representations of
125 protein motions with a convolutional neural network inputting pLM sequence embeddings
126 (**Fig. 1**). We considered the purely sequence-based pLM ESM2 (Lin et al., 2023) and two
127 structure-aware pLMs, namely ESM3 (Hayes et al., 2024) and ProstT5 (Heinzinger et al.,
128 2023). ESM3 is the largest model (**Table S1**), and it can condition on and reconstruct
129 several protein sequence and structural properties. ProstT5, the smallest model (**Table**
130 **S1**), is a fine-tuned version of the sequence-only model T5 that translates amino acid
131 sequences into sequences of discrete structural states and reciprocally. We trained and
132 tested SeaMoon on over $\sim 17\,000$ experimental conformational collections representing
133 a non-redundant set of the PDB at 80% sequence similarity. We used the principal
134 components extracted from these collections as ground-truth linear motions to which we
135 compared SeaMoon predicted 3D vectors. The latter are not anchored on a particular
136 conformation and may be in any arbitrary orientation. To allow for a fair comparison,
137 we determined the optimal rotation and scaling between the ground-truth and predicted
138 vectors before computing the error between them (see *Methods* for details). Based on
139 visual inspection, we considered predictions as acceptable when their normalised sum-
140 of-squares error (NSSE) was smaller than 0.6 (**Fig. 1**). See **Fig. S1** for illustrative
141 examples of different error levels. By comparison, random predictions typically display
142 errors above 0.9 (**Fig. S2**). SeaMoon is highly computationally efficient. It took 12s to
143 predict 3 motions for each of 1 121 test proteins on a desk computer equipped with Intel
144 Xeon W-2245 @ 3.90 GHz.

145 SeaMoon predicts motions from sequences across diverse protein 146 families

147 SeaMoon predicted at least one acceptable linear motion for each of 300 test proteins from
148 the purely sequence-based ESM2 embeddings (**Table I** and **Fig. 2A**). Its performance
149 was comparable to that of the purely geometry-based unsupervised NMA. SeaMoon suc-
150 cess rate improved by 25-40% when inputting structurally-informed embeddings com-
151 puted by ESM3 or ProstT5, outperforming the NMA by a large margin (**Table I** and
152 **Fig. 2A**). ProstT5, with the smallest number of parameters and embedding dimensions
153 (**Table S1**), yielded the best overall performance (**Fig. 2A**, paired Wilcoxon signed-
154 rank test p-values $< 10^{-6}$ and $< 10^{-9}$ with respect to ESM3 and ESM2, respectively). In
155 addition, we observed a boost in performance by up to 10% upon stimulating the model
156 to learn a one-sequence-to-many-motions mapping (**Table I** and **Fig. 2A**). More specifi-
157 cally, we augmented the training data by using multiple (up to 5) reference conformations
158 per experimental collection (**Table S2**). While the pLM embeddings within a collection
159 should be highly similar, the extracted motions may differ substantially from one refer-
160 ence to another (Lombard et al., 2024). The positive impact of this data augmentation
161 strategy was most visible for the ESM-based version of SeaMoon (**Table I** and **Fig. 2A**).

162 SeaMoon effectively generalised to unseen proteins across diverse families (**Table I**,
163 **Fig. 2B**, and **Fig. S4-5**). It produced high-quality predictions at different levels of
164 similarity to the training set, which we can interpret as varying difficulty levels. For
165 instance, SeaMoon-ESM2(x5) almost perfectly recapitulated the motions of antibodies
166 (**Fig. S5A**), a class of proteins well represented in both train and test sets. Beyond such
167 easy cases, SeaMoon-ESM2(x5) could transfer knowledge between proteins with similar

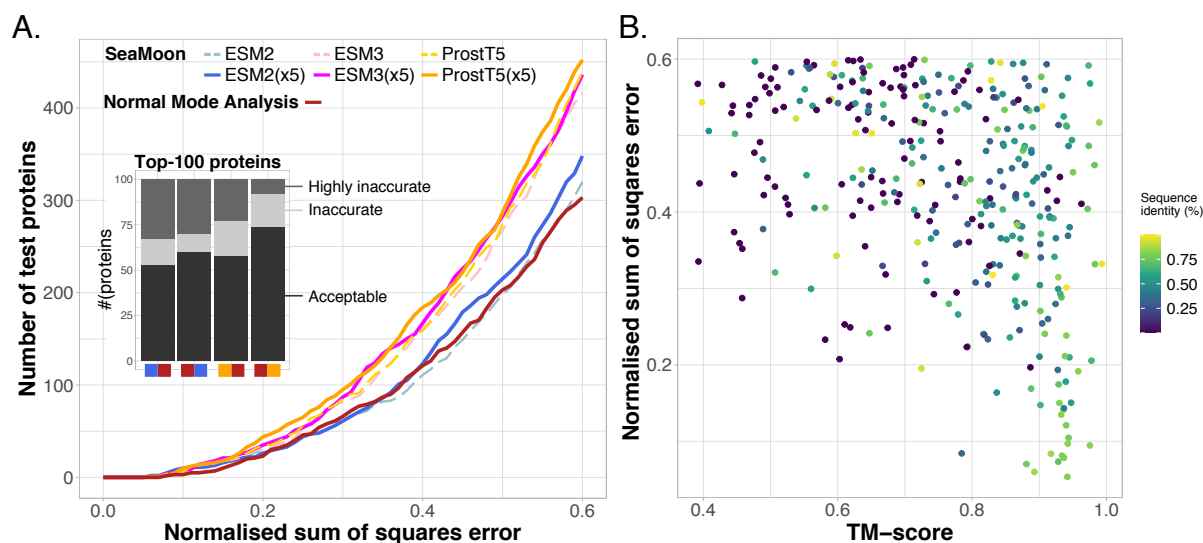


Figure 2: **SeaMoon performance and generalisation capability.** We report the NSSE of the best match between 3 predictions and 3 ground-truth motions for each of the 121 test proteins. **A.** Cumulative NSSE for six different versions of SeaMoon and for the NMA. We tested three pLMs, namely ESM2, ESM3 and ProstT5, and a data augmentation strategy with 5 training samples per experimental collection (x5). We cropped the plot at $NSSE = 0.6$ for ease of visualisation; see Fig. S3 for the full curves. Inset: Agreement between a selection of methods. For instance, the first bar stack gives the numbers of proteins for which the NMA (right red square) produced acceptable ($NSSE < 0.6$), inaccurate ($0.6 < NSSE < 0.75$) or highly inaccurate ($NSSE > 0.75$) predictions among the top-100 proteins best-predicted by SeaMoon-ESM2(x5) (left blue square). **B.** NSSE computed for SeaMoon-ESM2(x5) in function of sequence and structural similarity to the training set.

Table I: **Performance and dependence on the similarity to the training set**

Method	Protocol	Number of proteins w. acceptable predictions	Correlation w. TM-Score	Correlation w. sequence id.
SeaMoon	ESM2	320 (29%)	-0.35	-0.20
	ESM2(x5)	348 (31%)	-0.39	-0.26
	ESM3	416 (37%)	-0.31	-0.18
	ESM3(x5)	436 (39%)	-0.38	-0.22
	ProstT5	439 (39%)	-0.32	-0.12
	ProstT5(x5)	452 (40%)	-0.37	-0.20
NMA		303 (27%)	-0.09	0.03

We consider predictions as acceptable if their normalised sum-of-squares error is smaller than 0.6. The highest success rate is highlighted in bold.

168 3D folds but highly divergent sequences. The ATP-binding cassette (ABC) transporter
 169 superfamily provides an illustrative example of this intermediate difficulty (**Fig. S5B**).
 170 SeaMoon-ESM2(x5) accurately predicted the opening-closing motion of a putative ABC
 171 transporter from *Campylobacter jejuni* (**Fig. S5B**, 5T1PE, $NSSE = 0.33$) that does not
 172 have any detectable sequence similarity with the training set. This motion is character-

173 istic of the “Venus Fly-trap” mechanism for transporting sugars (Chandravanshi et al.,
174 2020) and is shared with a structurally similar ABC transporter from the training set
175 (**Fig. S5B**, 7C68B, TM-score = 0.83). At the most difficult level, SeaMoon-ESM2(x5)
176 successfully captured the motions of proteins completely unrelated to the training set,
177 such as the benzoyl-coenzyme A reductase from *Geobacter metallireducens* (**Fig. S5C**,
178 4Z3ZF, $NSSE = 0.37$).

179 SeaMoon complementary to the normal mode analysis

180 We investigated the extent of the agreement between the purely sequence-based version of
181 SeaMoon and the purely geometry-based NMA (**Fig. 2A**, inset, and **Fig. S6**). Among
182 the top-100 proteins best-predicted by SeaMoon-ESM2(x5), about half exhibit motions
183 accessible to the NMA (**Fig. 2A**, inset). Most of these motions involve a large portion of
184 the protein (median collectivity $\kappa = 0.69$) and correspond to large conformational changes
185 (median deviation of 5.1Å). They include functional opening-closing motions of virulence
186 factors, thermophilic proteins, metalloenzymes, periplasmic binding proteins, dehydroge-
187 nases, glutamate receptors, and antibodies (see **Fig. S7** for illustrative examples). On
188 the other hand, the NMA performed extremely poorly for a third of SeaMoon-ESM2(x5)
189 top-100 ($NSSE > 0.75$, see **Fig. 2A**, inset). The associated motions tend to be localised
190 with median collectivity $\kappa = 0.20$.

191 The bacterial toxins PemK and protective antigen (PA) from anthrax illustrate SeaMoon’s
192 capability to go beyond the NMA physics-based inference for highly localised motions
193 and fold-switching deformations (**Fig. 3**). SeaMoon-ESM2(x5) captured the PemK’s
194 loop L12 motion with high precision (**Fig. 3A**, $NSSE = 0.24$) whereas the NMA failed
195 to delineate the mobile region in the protein and to infer its direction of movement (**Fig.**
196 **3A**, in red). This highly localised motion ($\kappa = 0.17$) plays a decisive role in regulating
197 PemK RNase activity by promoting the formation of the PemK-PemI toxin-antitoxin
198 (Kim et al., 2022). In the anthrax protective antigen, SeaMoon-ESM2(x5) accurately
199 predicted the relative motion amplitudes and directions of an 80 residue-long region that
200 detaches from the rest of the protein upon forming an heptameric pore **Fig. 3B**). By
201 contrast, the NMA predicted a breathing motion poorly approximating the ground-truth
202 one (**Fig. 3B**), likely due to its assumption that proteins behave as elastic networks.
203 PA’s $\sim 30\text{Å}$ -large conformational transition is essential for the translocation of the bac-
204 terium’s edema and lethal factors to the host cell (Machen et al., 2021). PemK and PA do
205 not have any detectable sequence similarity to the training set. SeaMoon likely leveraged
206 information coming from training proteins with similar folds and functions from other
207 bacteria (Anderson et al., 2020; Dhanasingh et al., 2021).

208 Reciprocally, SeaMoon covered 60% of the top-100 proteins best-predicted by the
209 NMA with ESM2 embeddings, and up to 75% with ProstT5 embeddings (**Fig. 2A**,
210 inset, and **Fig. S6**). Using implicit structural knowledge allowed recovering elastic
211 motions such as that exhibited by the mammalian plexin A4 ectodomain (**Fig. S8**,
212 $NSSE = 0.28$). Taken together, SeaMoon-ProstT5(x5) and the NMA approximated the
213 motions of 554 test proteins (out of 1121, 49%) with reasonable accuracy (**Table. I**). This
214 result suggests that combining SeaMoon transfer learning approach with the physics- and
215 geometry-based NMA could be a valuable strategy.

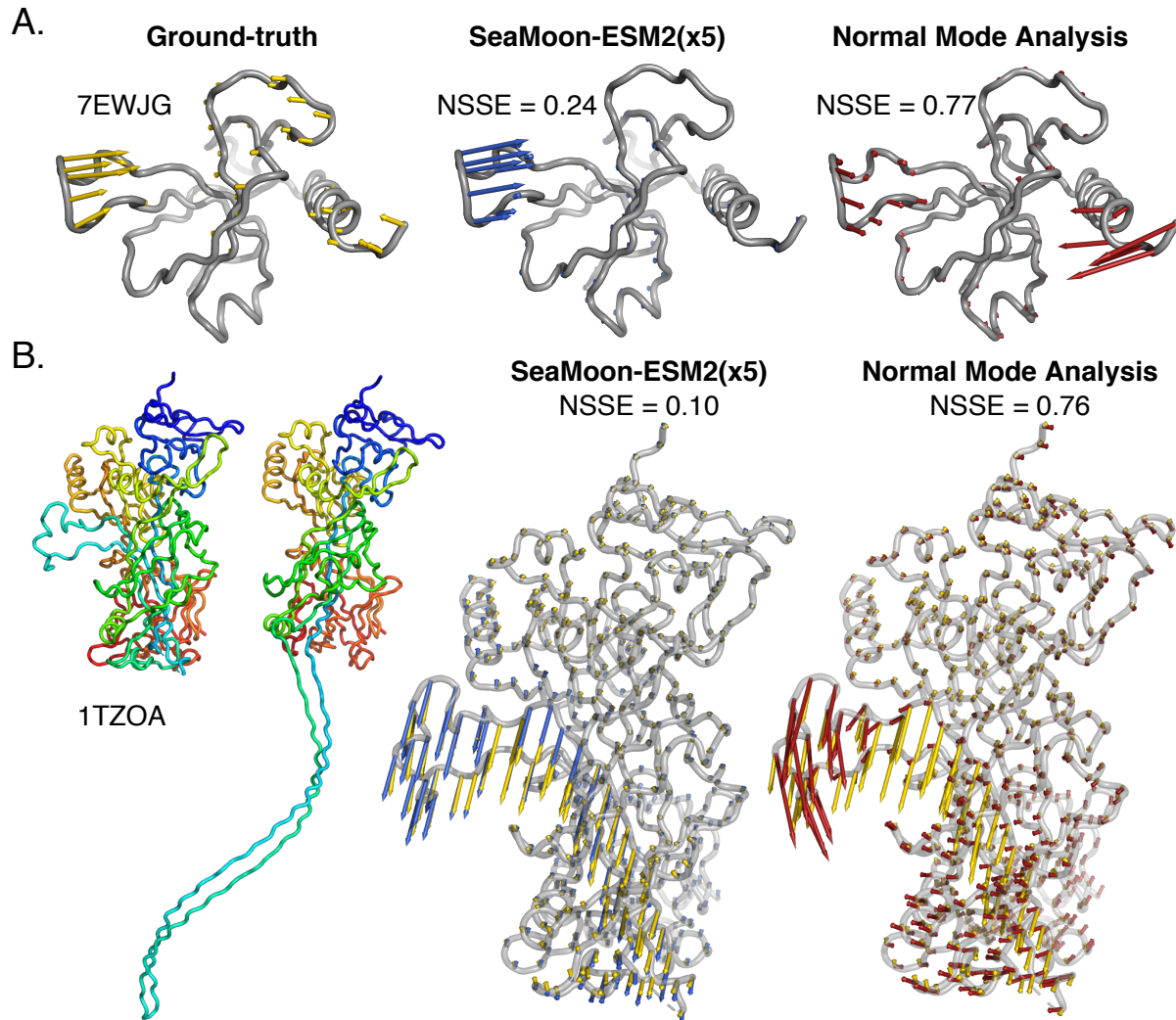


Figure 3: **Examples of motions well predicted by SeaMoon and not by the NMA.** The arrows depicted in yellow, blue and red on the grey 3D structures represent the ground-truth motions and the best-matching predictions from SeaMoon-ESM2(x5) and the NMA, respectively. **A.** Bacterial toxin PemK (PDB code: 7EWJ, chain G) from the test set. It does not have any detectable sequence similarity to the training set **B.** Anthrax protective antigen (PDB code: 1TZO, chain A) from the validation set. We show the two most extreme conformations of the collection on the left, colored according to the residue index, from the N-terminus in blue, to the C-terminus in red. The closest homolog from the training set shares 35% sequence similarity.

216 SeaMoon can recapitulate entire motion subspaces

217 Beyond assessing individual predictions, we evaluated the global similarities between
218 predicted and ground-truth 3-motion subspaces focusing on the test proteins for which
219 SeaMoon produced at least one acceptable prediction (**Table I**). We found that SeaMoon
220 motion subspaces were fairly similar to the ground-truth ones, with a Root Mean Square
221 Inner Product (RMSIP) (Amadei et al., 1999; Leo-Macias et al., 2005; David and Ja-
222 cobs, 2011) higher than 0.5, for almost two thirds of these proteins. We observed an
223 excellent correspondence for a dozen proteins, *e.g.*, the *Mycobacterium* phage Ogopogo
224 major capsid protein (**Fig. 4** and **Fig. S9**). The purely sequence-based SeaMoon-

225 ESM2(x5) achieved an RMSIP of 0.75 on this protein, and the structure-aware SeaMoon-
226 ProstT5(x5) reached 0.82. SeaMoon-ProstT5(x5) first, second and third predicted mo-
227 tions had a Pearson correlation of 0.93, 0.73 and 0.75 with the first, third and second
228 ground-truth principal components, respectively (**Fig. 4A**). The associated NSSE were
229 all smaller than 0.5 (**Fig. 4B**). By inspecting the training set, we could identify sev-
230 eral major capsid proteins from other bacteriophages sharing the same HK97-like fold
231 as the Ogotogo one (TM-score up to 0.78), despite relatively low sequence similarity
232 (up to 34%). The ability of SeaMoon to recapitulate the Ogotogo protein entire motion
233 subspace with reasonable accuracy likely reflects the high conservation of major capsid
234 protein dynamics upon forming icosahedral shells (Podgorski et al., 2023).

235 Contributions of the inputs and design choices

236 We investigated the contribution of SeaMoon inputs, architecture and objective function
237 to its success rate through an ablation study, starting from SeaMoon-ProstT5 baseline
238 model (**Table S3** and **Fig. S10**). Inputting random matrices instead of pre-trained pLM
239 embeddings or using only positional encoding had the most drastic impacts. Still, we ob-
240 served that the network can produce accurate predictions for over 100 proteins in this
241 extreme situation (**Fig. S10**, in grey). Annihilating sequence embedding context by set-
242 ting all convolutional filter sizes to 1 also had a dramatic impact, reducing to success rate
243 from 40 to 25% (**Table S3** and **Fig. S10**). Moreover, a 7-layer transformer architecture
244 (see *Methods*) underperformed SeaMoon’s convolutional neural network, despite having
245 roughly the same number of free parameters (**Fig. S10**, in brown). Finally, disabling
246 either sign flip or reflection (*i.e.*, pseudo-rotation) or permutation when computing the
247 loss degraded the performance by 6 to 15% (**Fig. S10**, in light green). This result un-
248 derlines the utility of implementing a permissive and flexible comparison of the predicted
249 and ground-truth motions during training.

250 SeaMoon practical utility to deform protein structures

251 SeaMoon does not use any explicit 3D structural information during inference. Its pre-
252 dictions are independent of the global orientation of any protein conformation, making it
253 impractical to directly use them to deform protein structures. To partially overcome this
254 limitation, we propose an unsupervised procedure to orient SeaMoon predicted vectors
255 with respect to a given protein 3D conformation. This method exploits the rotational
256 constraints of the ground-truth principal components. Namely, the total angular velocity
257 of the reference conformation subjected to a ground-truth principal component is zero (see
258 *Methods*). Therefore, we determine the rotation that must be applied to the predicted
259 motion vectors to minimize the total angular velocity of a target conformation.

260 This strategy proved successful for the vast majority of SeaMoon’s highly accurate
261 predictions. SeaMoon-ProstT5(x5) predicted motion vectors, oriented to minimise an-
262 gular velocity, exhibit an acceptable error (< 0.6) in 85% of cases where the optimal
263 alignment with the ground truth results in $NSSE < 0.3$. This result indicates that pre-
264 dictions that approximate well the ground-truth principal components also preserve their
265 properties. The human ABC transporter sub-family B member 6 gives an illustrative
266 example where the third predicted motion vector approximates the first ground-truth
267 principal component with $NSSE = 0.20$ upon optimal alignment and 0.22 upon angular
268 velocity minimisation (**Fig. 4C-E**). Overall, the procedure allowed for correctly orienting

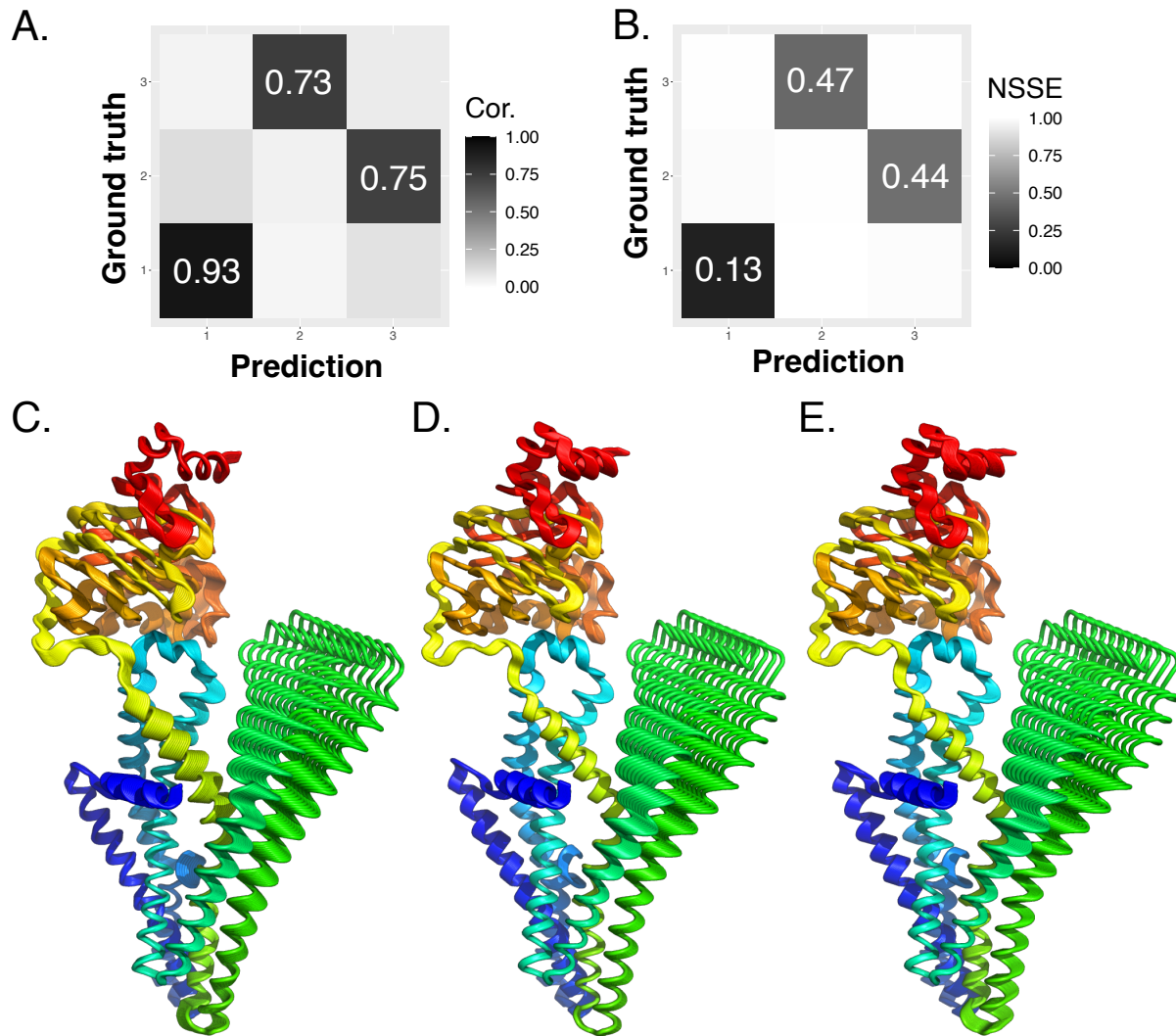


Figure 4: **Motion subspace comparison and deformation trajectories.** **A-B.** Ogopogo major capsid protein motion subspace. PDB code: 8ECN, chain B. **A.** Pairwise similarities measured as Pearson correlations between the ground-truth motions and SeaMoon-ProstT5(x5) predictions. **B.** Pairwise discrepancies measured as NSSE. **C-E.** Trajectories of a human ABC transporter (PDB code: 7D7R, chain A) deformed along its first ground-truth principal component (A) and the best-matching SeaMoon-ProstT5(x5) prediction (B-C). **B.** The prediction is optimally aligned with the ground truth. **C.** The orientation of the prediction minimises the protein conformation's angular velocity. Each trajectory comprises 10 conformations coloured from blue at the N-terminus to red at the C-terminus.

269 acceptable predictions for 215 test proteins.

270 Note that this post-processing increases computing time significantly, from 12s to 24m
271 over the 1 121 test proteins on a desk computer equipped with Intel Xeon W-2245 @ 3.90
272 GHz.

273 Methods

274 Datasets

275 To generate training data, we constructed a non-redundant set of conformational col-
276 lections representing the whole PDB (as of June 2023) using DANCE (Lombard et al.,
277 2024). To ensure high quality of the data, we replaced the raw PDB coordinates with
278 their updated and optimised versions from PDB-REDO whenever possible (Joosten et al.,
279 2014). We used a stringent setup where each conformational collection is specific to a set
280 of close homologs. Specifically, any two protein chains belonging to the same collection
281 share at least 80% sequence identity and coverage. We filtered out the collections with
282 too few or too many data points. Namely, we asked for at least 4 and at most 500 con-
283 formations and a representative protein chain comprising between 30 and 1 000 residues.
284 We further retained only C α atoms (option -c) and used coordinate weights to account
285 for uncertainty (option -w).

286 For each collection, DANCE extracted the $K = 3$ principal components contribut-
287 ing the most to its total positional variance (Lombard et al., 2024). We interpret these
288 components as the main linear motions explaining the collection’s conformational di-
289 versity. Namely, the k th principal component defines a set of 3D displacement vectors
290 $\{\bar{x}_{ik}^{\text{GT}}, i = 1, 2, \dots, L\}$ for the L protein residues’ C α atoms. We normalised these vectors
291 to facilitate their comparison across different proteins, such that $\sum_{i=1}^L \|\bar{x}_{ik}^{\text{GT}}\|^2 = L$. We
292 further applied three filtering criteria with the aim of excluding collections with low di-
293 versity or highly non-linear complex deformations: (i) maximum Root Mean Squared
294 Deviation (RMSD) between any two conformations of at least 2 Å, (ii) first principal
295 component (main linear motion) contributing at least 80% of the total variance and (iii)
296 involving at least 12 residues, *i.e.*, $L \times \kappa \geq 12$, where κ is the collectivity of the principal
297 component (see definition below). This operation resulted in 7 339 collections, randomly
298 split between train (70%), validation (15%) and test (15%) sets.

299 DANCE makes use of a reference conformation to superimpose the C α atoms’ 3D
300 coordinates and centre them prior to extracting motions with PCA. By default, the refer-
301 ence corresponds to the protein chain with the most representative amino acid sequence
302 (Lombard et al., 2024). In order to augment the data, we defined up to 4 alternative
303 reference conformations, in addition to the default one (option -n 5). At each iteration,
304 DANCE chose the new reference conformation as the one displaying the highest RMSD
305 from the previous one. This strategy maximises the impact of changing the reference and
306 thus the diversity of the extracted motions.

307 Model Specifications

308 Input features

309 SeaMoon takes as input embeddings computed from pre-trained pLMs, namely Evolution-
310 ary Scale Models ESM2-T33-650M-UR50 (Lin et al., 2023) and ESM3-small (1.4B) (Hayes
311 et al., 2024), as well as Protein sequence-structure T5 (Heinzinger et al., 2023). ESM2-
312 T33-650M-UR50 is a BERT (Devlin et al., 2018) style 650-million-parameter encoder-only
313 transformer architecture trained on all clusters from Uniref50 (Suzek et al., 2015, 2007),
314 a version of UniProt (Consortium, 2022) clustered at 50% sequence similarity, augmented
315 by sampling sequences from the Uniref90 clusters of the representative chains (excluding
316 artificial sequences). ESM3-small (1.4B) is a transformer-based (Vaswani et al., 2017) all-

317 to-all generative architecture that both conditions on and generates a variety of different
318 tracks representing protein sequence, secondary and tertiary structure, solvent accessibil-
319 ity and function. It was trained on over 2.5 billion natural proteins collected from sequence
320 and structure databases, including UniRef, MGnify (Richardson et al., 2023), OAS (Olsen
321 et al., 2022) and the PDB (Berman et al., 2000), augmented with synthetic sequences
322 generated by an inverse folding model (Hayes et al., 2024). Protein sequence-structure
323 T5 is a bilingual pLM trained on a high-quality clustered version of the AlphaFold Pro-
324 tein Structure Database (Barrio-Hernandez et al., 2023; Varadi et al., 2021) to translate
325 1D sequences of amino acids into 1D sequences of 3Di tokens representing 3D structural
326 states (Van Kempen et al., 2024) and vice versa. The 3Di alphabet, introduced by the
327 3D-alignment method Foldseek (Van Kempen et al., 2024), describes tertiary contacts be-
328 tween protein residues and their nearest neighbours. This 1D discretised representation
329 of 3D structures is sensitive to fold change but robust to conformational rearrangements.
330 Protein sequence-structure T5 expands on ProtT5-XL-U50 (Elnaggar et al., 2022), an
331 encoder-decoder transformer architecture (Raffel et al., 2020) trained on reconstructing
332 corrupted amino acids from the Big Fantastic Database (Steinegger et al., 2019) and
333 UniRef50. Throughout the text, we refer to these pLMs as ESM2, ESM3 and ProstT5,
334 respectively. We used the pre-trained pLMs as is, without fine-tuning their weights, and
335 we gave them only amino acid sequences as input.

336 **Model’s architecture**

337 SeaMoon’s architecture is a convolutional neural network (LeCun et al., 2015) taking as
338 input a sequence embedding of dimensions $L \times d$, with L the number of protein residues
339 and d the representation dimension of the chosen pLM, namely 1 280 for ESM2, 1 536 for
340 ESM3, and 1 024 for ProstT5, and outputting K predicted tensors of dimensions $L \times 3$. It
341 comprises a linear layer followed by two hidden 1-dimensional convolutional layers with
342 filter sizes of 15 and 31, respectively, and finally K parallel linear layers (**Table S1**).
343 SeaMoon’s convolutional architecture allows handling sequences of any arbitrary length
344 L and preserving this dimension throughout the network. All layers were linked through
345 the LeakyReLU activation function (Maas et al., 2013), as well as 80% dropout (Srivastava
346 et al., 2014). We experimented with other types of architectures, including those based
347 on sequence transformers, and chose the one based on CNNs as it demonstrated the
348 maximum accuracy at a reasonable number of trained parameters. Please see **Table S3**
349 and **Fig. S10** for more details. We implemented the models in PyTorch (Paszke et al.,
350 2019) v2.1.0 using Python 3.11.9.

351 By design, the SeaMoon model predicts the K motion tensors in a latent space that is
352 invariant to the protein’s actual 3D orientation. To align these predictions with a given
353 3D conformation, additional information, such as the ground-truth motions, is required,
354 as explained below.

355 **Loss function**

356 We aim to minimise the discrepancy between the predicted tensor X and the ground-truth
357 tensor X^{GT} , both of dimensions $L \times K \times 3$, expressed as a weighted aligned sum-of-squares
358 error loss,

$$\mathcal{L} = \frac{1}{L} \min_{R,S,P} \left(\sum_{i=1}^L w_i \|R(PX_i^{\text{GT}})^T - (SX_i)^T\|_F^2 \right), \quad (1)$$

359 where X_i defines the set of K 3D displacements vectors $\{\vec{x}_{ik} \equiv (X_{i,k,\cdot})^T, k = 1, 2, \dots, K\}$
360 predicted for the C α atom of residue i , X_i^{GT} defines the corresponding ground-truth
361 3D displacement vector set, $\|\cdot\|_F$ designates the Frobenius norm, and w_i is a weight
362 reflecting the confidence in the ground-truth data for residue i (Lombard et al., 2024). It is
363 computed as the proportion of conformations in the experimental collection with resolved
364 3D coordinates for residue i . The matrices R , of dimension 3×3 , and P , of dimension
365 $K \times K$, allow for rotating and permuting the ground-truth vectors to optimally align
366 them with the predicted ones. We chose to apply the transformations to the ground-
367 truth vectors for gradient stability. We allow for rotations R because SeaMoon relies
368 solely on a protein sequence embedding as input. Its predictions are not anchored in a
369 particular 3D structure and hence, they may be in any arbitrary orientation. We allow for
370 permutation P to stimulate knowledge transfer across conformational collections. The
371 rationale is that a motion may be shared between two collections without necessarily
372 contributing to their positional variance to the same extent. Additionally, we allow for
373 scaling predictions with the $K \times K$ diagonal matrix S , so that SeaMoon can focus on
374 predicting only the relative motion amplitudes between the amino acid residues.

375 In practice, we first jointly determine the optimal permutation P and rotation R
376 of the ground-truth 3D vectors. We test all possible permutations, and, for each, we
377 determine the best rotation by solving the orthogonal Procrustes problem (Gower and
378 Dijkstra, 2004; Schönemann, 1966). We shall note that the optimal solution may be
379 a pseudo-rotation, *i.e.*, $\det(R) = -1$, which corresponds to the combination of a rotation
380 and an inversion. The loss can then be reformulated as,

$$\mathcal{L} = \frac{1}{L} \min_S \left(\sum_{k=1}^K \sum_{i=1}^L w_i \|\vec{x}_{ik}^{\text{GT-trans}} - S_{kk} \vec{x}_{ik}\|^2 \right), \quad (2)$$

381 where $\vec{x}_{ik}^{\text{GT-trans}}$ is the ground-truth 3D displacement vector for residue i matching the
382 predicted 3D vector \vec{x}_{ik} and aligned with it, and $S_{kk} \in \mathbb{R}$ is the k th scaling coefficient,
383 *i.e.* the k th non-null term of the diagonal scaling matrix S . The optimal value for S_{kk} is
384 computed as,

$$S_{kk} = \frac{\sum_{i=1}^L w_i (\vec{x}_{ik}^{\text{GT-trans}})^T \vec{x}_{ik}}{\sum_{i=1}^L w_i \|\vec{x}_{ik}\|^2}. \quad (3)$$

385 Training

386 We trained six models (**Table S2**) to predict $K = 3$ motions using the Adam optimizer
387 (Kingma and Ba, 2014) with a learning rate of 1e-02. We used a batch size of 64 input
388 sequences and employed padding to accommodate sequences of variable sizes in the same
389 batch. We trained for 500 epochs and kept the best model according to the performance
390 on the validation set.

391 Inference

392 We provide an unsupervised procedure to orient SeaMoon's predicted motions with re-
393 spect to a target 3D conformation \vec{C}_i during inference. This approach relies on the as-
394 sumption that correct predictions comply with the same rotational constraints as ground-
395 truth motions (see *Supplementary Methods*). Specifically, these constraints state that the
396 cross products between the positional 3D vectors of the reference conformation C^0 and
397 the 3D displacement vectors defined by a ground-truth principal component X_k^{GT} result

398 in a null vector,

$$\sum_{i=1}^L \vec{C}_i^0 \times \vec{x}_{ik}^{\text{GT}} = \vec{0}. \quad (4)$$

399 Assuming that the motion tensor X_k predicted by SeaMoon preserves this property, we
400 determine the rotation R that minimises the following cross-product,

$$\sum_{i=1}^L \vec{C}_i \times R\vec{x}_{ik} = \vec{0}. \quad (5)$$

401 This problem has at most four solutions and we solve it exactly using the symbolic
402 *wolframclient* package in Python. See *Supplementary Methods* for a detailed explanation.
403 In practice, we observe that these four solutions reduce to two pairs of highly similar
404 rotations.

405 Evaluation

406 We assessed SeaMoon predictions on each test protein from two different perspectives.
407 In the first assessment, we considered all $K \times K$ pairs of predicted and ground-truth
408 motions and estimated the discrepancy between the two motions within each pair after
409 optimally rotating and scaling them. We focused on the best matching pair for computing
410 success rates and illustrating the results. In the second assessment, we considered the
411 predicted and ground-truth motion subspaces at once and estimated their permutation-,
412 rotation- and scaling-invariant global similarity. In addition, we estimated discrepancies
413 and similarities between individual predicted and ground-truth motions after globally
414 matching and aligning the subspaces. We detail our evaluation metrics and procedures
415 in the following.

416 Normalised sum-of-squares error

At inference time, we estimate the discrepancy between the k th predicted motion and the
 l th ground-truth principal component by computing their weighted sum-of-squares error
under optimal rotation R^{opt} and scaling s^{opt} ,

$$SSE = \frac{1}{L} \sum_{i=1}^L w_i \|\vec{x}_{il}^{\text{GT-trans}} - s^{\text{opt}} \vec{x}_{ik}\|^2, \quad (6)$$

$$\text{with } \vec{x}_{il}^{\text{GT-trans}} = R^{\text{opt}} \vec{x}_{il}^{\text{GT}} \quad (7)$$

417 In the best-case scenario, the prediction is colinear to the transformed ground-truth,
418 $\vec{x}_{il}^{\text{GT-trans}} = c \vec{x}_{ik}$, $c \in \mathbb{R}$, such that $(\vec{x}_{il}^{\text{GT-trans}})^T \vec{x}_{ik} = \|\vec{x}_{il}^{\text{GT-trans}}\| \|\vec{x}_{ik}\| = c \|\vec{x}_{ik}\|^2$, $\forall i \in$
419 $1, 2, \dots, L$. By virtue of 3, the scaling coefficient s^{opt} will be equal to c , and thus, the error
420 will be null,

$$SSE_{\text{min}} = \frac{1}{L} \sum_{i=1}^L w_i \|\vec{x}_{il}^{\text{GT-trans}} - c \vec{x}_{ik}\|^2 = \frac{1}{L} \sum_{i=1}^L w_i \|c \vec{x}_{ik} - c \vec{x}_{ik}\|^2 = 0. \quad (8)$$

421 In the worst-case scenario, the prediction is orthogonal to the ground truth, such that
422 $(\vec{x}_{il}^{\text{GT-trans}})^T \vec{x}_{ik} = 0$, $\forall i \in 1, 2, \dots, L$. The scaling coefficient will be null and, hence, this
423 situation is equivalent to having a null prediction,

$$SSE_{\text{max}} = \frac{1}{L} \sum_{i=1}^L w_i \|\vec{x}_{il}^{\text{GT-trans}} - \vec{0}\|^2 = \frac{1}{L} \sum_{i=1}^L w_i \|\vec{x}_{il}^{\text{GT-trans}}\|^2. \quad (9)$$

424 The value of the raw error depends on the uncertainty of the ground-truth data. If all
425 conformations in the collection have resolved 3D coordinates for all protein residues, then
426 $w_i = 1, \forall i = 1, 2, \dots, L$ and the maximum error is $SSE_{max} = \frac{1}{L} \sum_{i=1}^L \|\vec{x}_{il}^{GT-trans}\|^2 = \frac{L}{L} = 1$.
427 As uncertainty in the ground-truth data increases, the associated errors will become
428 smaller. To ensure a fair assessment of the predictions across proteins, we normalise the
429 raw errors,

$$NSSE = \frac{SSE}{SSE_{max}}. \quad (10)$$

430 Estimation of sum-of-squares errors for random vectors

431 To compare SeaMoon results with a random baseline, we selected 14 ground-truth prin-
432 cipal components from the test set. We focused on proteins with maximum confidence,
433 *i.e.*, for which $w_i = 1, \forall i = 1, 2, \dots, L$. We started with a set of 10 components chosen
434 randomly. We then added the most localised component (collectivity $\kappa = 0.06$), the most
435 collective one ($\kappa = 0.85$), a component from the smallest protein (33 residues), and a com-
436 ponent from the longest one (662 residues). We generated 1000 random predictions for
437 each ground truth component and computed their sum-of-squares errors under optimal
438 rotation and scaling.

439 Subspace comparison

440 We estimated the similarity between the $K \times 3$ subspaces spanned by SeaMoon predictions
441 and the ground-truth principal components as their Root Mean Square Inner Product
442 (RMSIP) (Amadei et al., 1999; Leo-Macias et al., 2005; David and Jacobs, 2011). It
443 is computed as an average of the normalised inner products of all the vectors in both
444 subspaces,

$$\text{RMSIP} = \left(\frac{1}{K} \sum_{k=1}^K \sum_{l=1}^K \sum_{i=1}^L \frac{(\vec{x}_{ik}^{GT})^T \vec{x}_{il}^{\text{ortho}}}{\|\vec{x}_{ik}^{GT}\| \|\vec{x}_{il}^{\text{ortho}}\|} \right), \quad (11)$$

445 where $\vec{x}_{il}^{\text{ortho}}$ is obtained by orthogonalising SeaMoon predictions using the Gram–Schmidt
446 process. This operation ensures that the RMSIP ranges from zero for mutually orthog-
447 onalising subspaces to one for identical subspaces and avoids artificially inflating the
448 RMSIP due to redundancy in the predicted motions. We should stress that in practice,
449 this redundancy is limited and the motions predicted for a given protein never collapse
450 (**Fig. S11**). A RMSIP score of 0.70 is considered an excellent correspondence while a
451 score of 0.50 is considered fair (Amadei et al., 1999).

452 While the RMSIP is invariant to permutations and rotations, the individual inner
453 products, reflecting similarities between pairs of motions, are not. For interpretability
454 purposes, we maximised these pairwise similarities through the following procedure:

- 455 1. compute the NSSE for all pairs of predictions and ground-truth principal compo-
456 nents, under optimal rotation and scaling, as in 7,
- 457 2. orthogonalise the predictions in the order of their losses, from the best-matching
458 prediction to the worst-matching one,
- 459 3. determine the optimal global rotation of the ordered set of matching ground-truth
460 components onto the ordered set of orthogonalised predictions,

461 4. compute all pairwise normalised inner products and the corresponding RMSIP, and
462 all pairwise NSSE under optimal scaling.

463 Comparison with the normal mode analysis

464 We compared SeaMoon performance with the physics-based unsupervised normal mode
465 analysis (NMA) (Hayward and Go, 1995). The NMA takes as input a protein 3D structure
466 and builds an elastic network model where the nodes represent the atoms and the edges
467 represent springs linking atoms located close to each other in 3D space. The normal modes
468 are obtained by diagonalizing the mass-weighted Hessian matrix of the potential energy
469 of this network. We used the highly efficient NOLB method (Hoffmann and Grudinin,
470 2017) to extract the first $K = 3$ normal modes from the test protein 3D conformations.
471 We retained only the $C\alpha$ atoms, as for the principal component analysis, and defined
472 the edges in the elastic network using a distance cutoff of 10\AA . We enhanced the elastic
473 network dynamical potential by excluding edges corresponding to small contact areas
474 between protein segments. We detected them as disconnected patches in the contact
475 map using HOPMA (Laine and Grudinin, 2021). Contrary to SeaMoon predictions, the
476 orientation of the NMA predictions is not arbitrary and thus, we do not need to align the
477 ground-truth components onto them.

478 Motion properties

479 Contribution

480 We estimate the contribution of the $L \times 3$ ground-truth principal component X_k^{GT} to the
481 total positional variance as its normalised eigenvalue, $\frac{\lambda_k}{\sum_l \lambda_l}$.

482 Collectivity

483 We estimate the collectivity (Brüschweiler, 1995; Tama and Sanejouand, 2001) of the
484 $L \times 3$ predicted or ground-truth motion tensor X_k as,

$$\kappa(X_k) = \frac{1}{L} \exp \left(- \sum_{i=1}^L \sum_{j=1}^3 X_{ijk}^2 \log X_{ijk}^2 \right), \quad (12)$$

485 with L the number of residues. If $\kappa(\mathbf{v}) = 1$, then the corresponding motion is maximally
486 collective and has all the atomic displacements identical. In case of an extremely localised
487 motion, where only one single atom is affected, the collectivity is minimal and equals to
488 $1/L$.

489 Conclusion

490 This proof-of-concept study explores the extent to which protein sequences encode func-
491 tional motions. SeaMoon reconstructs these motions within an invariant subspace directly
492 from sequence-based pLM embeddings. Our results indicate that incorporating structure-
493 aware input embeddings significantly improves the success rate. Moreover, they highlight
494 SeaMoon’s ability to transfer knowledge about motions across distant homologs, lever-
495 aging the universal representation space of pLMs. However, the framework’s capacity to
496 predict entirely novel motions has yet to be fully assessed.

497 SeaMoon’s transfer learning approach complements unsupervised methods that rely
498 solely on the 3D geometry of protein structures, such as Normal Mode Analysis (NMA).
499 Future work will focus on integrating these two sources of information into a unified,
500 end-to-end framework. Incorporating explicit structural information for a target protein
501 could resolve the ambiguity in orienting predicted motions without requiring ground-truth
502 knowledge.

503 One current limitation is the scarcity of functional motions in the training set, raising
504 concerns about its accuracy and completeness. Both SeaMoon and NMA struggle to
505 predict certain motions, suggesting that these may lack biological or physical relevance.
506 Conversely, SeaMoon could be used to assess the evolutionary conservation of motions.
507 Another limitation of the current approach is its reliance on a linear description of protein
508 motion subspaces. Linear principal components are insufficient for describing complex
509 loop deformations or large rearrangements of secondary structures. Introducing non-
510 linearity could yield more realistic motion predictions. Future work will address these
511 issues, potentially augmenting the training set with *in silico* generated data, such as
512 motions derived from MD and NMA simulations, or protein conformations predicted by
513 AlphaFold.

514 Despite these limitations, the current findings offer valuable insights for integrative
515 structural biology. SeaMoon provides a compact representation of continuous structural
516 heterogeneity in proteins, enabling the sampling of conformations through a generative
517 model. Additionally, the estimated motion subspaces can be used to compute protein
518 conformational entropy. Lastly, our framework is highly versatile, featuring a lightweight,
519 trainable deep learning architecture that does not depend on fine-tuning a large pre-
520 trained model. This flexibility allows users to easily adapt the system to new input pLM
521 embeddings without modifying the model architecture.

522 Declaration

523 Data and code availability

524 The source code and model weights of this work are freely available at <https://github.com/PhyloSoFS-Team//seamoon>. The data used for development of SeaMoon and SeaMoon
525 predictions are freely available at Zenodo.
526

527 Acknowledgments

528 The Sorbonne Center for Artificial Intelligence (SCAI) provided a salary to VL and
529 computational resources. This work has also been co-funded by the European Union
530 (ERC, PROMISE, 101087830). Views and opinions expressed are however those of the
531 author(s) only and do not necessarily reflect those of the European Union or the European
532 Research Council. Neither the European Union nor the granting authority can be held
533 responsible for them.

534 Author contributions

535 S.G. and E.L. designed research and supervised the project. V.L. designed the model’s
536 architecture and carried out its implementation. S.G. and D.T. wrote the proofs and
537 problem formalisation for orienting predictions with respect to a protein conformation

538 with feedback from E.L.. D.T. implemented the solver. V.L., E.L. and S.G. produced
539 and analysed the results. E.L. wrote the manuscript with input, support and feedback
540 from all authors. All authors edited, read, and approved the final manuscript.

541 Competing interests

542 The author(s) declare no competing interests.

543 References

- 544 Lane, T. J. Protein structure prediction has reached the single-structure frontier. *Nature*
545 *Methods* **2023**, *20*, 170–173.
- 546 Miller, M. D.; Phillips, G. N. Moving beyond static snapshots: Protein dynamics and the
547 Protein Data Bank. *Journal of Biological Chemistry* **2021**, *296*.
- 548 Henzler-Wildman, K.; Kern, D. Dynamic personalities of proteins. *Nature* **2007**, *450*,
549 964–972.
- 550 Jumper, J. et al. Highly accurate protein structure prediction with AlphaFold. *Nature*
551 **2021**, *596*, 583–589.
- 552 Sala, D.; Engelberger, F.; Mchaourab, H.; Meiler, J. Modeling conformational states of
553 proteins with AlphaFold. *Current Opinion in Structural Biology* **2023**, *81*, 102645.
- 554 Faezov, B.; Dunbrack Jr, R. L. AlphaFold2 models of the active form of all 437
555 catalytically-competent typical human kinase domains. *bioRxiv* **2023**, 2023–07.
- 556 Heo, L.; Feig, M. Multi-state modeling of G-protein coupled receptors at experimental
557 accuracy. *Proteins: Structure, Function, and Bioinformatics* **2022**, *90*, 1873–1885.
- 558 Wallner, B. AFsample: improving multimer prediction with AlphaFold using massive
559 sampling. *Bioinformatics* **2023**, *39*, btad573.
- 560 Wallner, B. Improved multimer prediction using massive sampling with AlphaFold in
561 CASP15. *Proteins: Structure, Function, and Bioinformatics* **2023**, *91*, 1734–1746.
- 562 Johansson-Åkhe, I.; Wallner, B. Improving peptide-protein docking with AlphaFold-
563 Multimer using forced sampling. *Frontiers in Bioinformatics* **2022**, *2*, 85.
- 564 others,, et al. Impact of AlphaFold on structure prediction of protein complexes: The
565 CASP15-CAPRI experiment. *Proteins: Structure, Function, and Bioinformatics* **2023**,
566 *91*, 1658–1683.
- 567 Brysbaert, G.; Raouraoua, N.; Mirabello, C.; Véry, T.; Blanchet, C.; Wallner, B.;
568 Lensink, M. MassiveFold: unveiling AlphaFold’s hidden potential with optimized and
569 parallelized massive sampling. **2024**,
- 570 Sfriso, P.; Duran-Frigola, M.; Mosca, R.; Emperador, A.; Aloy, P.; Orozco, M. Residues
571 coevolution guides the systematic identification of alternative functional conformations
572 in proteins. *Structure* **2016**, *24*, 116–126.

- 573 Benner, S. A.; Gerloff, D. Patterns of divergence in homologous proteins as indicators of
574 secondary and tertiary structure: a prediction of the structure of the catalytic domain
575 of protein kinases. *Advances in Enzyme Regulation* **1991**, *31*, 121–181.
- 576 Göbel, U.; Sander, C.; Schneider, R.; Valencia, A. Correlated mutations and residue
577 contacts in proteins. *Proteins: Structure, Function, and Bioinformatics* **1994**, *18*, 309–
578 317.
- 579 Ortiz, A. R.; Kolinski, A.; Rotkiewicz, P.; Ilkowski, B.; Skolnick, J. Ab initio folding of
580 proteins using restraints derived from evolutionary information. *Proteins: Structure,
581 Function, and Bioinformatics* **1999**, *37*, 177–185.
- 582 Lapedes, A. S.; Giraud, B. G.; Liu, L.; Stormo, G. D. Correlated mutations in models of
583 protein sequences: phylogenetic and structural effects. *Lecture Notes-Monograph Series*
584 **1999**, 236–256.
- 585 Giraud, B.; Heumann, J. M.; Lapedes, A. S. Superadditive correlation. *Physical Review
586 E* **1999**, *59*, 4983.
- 587 Thomas, J.; Ramakrishnan, N.; Bailey-Kellogg, C. Graphical models of residue coupling
588 in protein families. Proceedings of the 5th international workshop on Bioinformatics.
589 2005; pp 12–20.
- 590 Weigt, M.; White, R. A.; Szurmant, H.; Hoch, J. A.; Hwa, T. Identification of direct
591 residue contacts in protein–protein interaction by message passing. *Proceedings of the
592 National Academy of Sciences* **2009**, *106*, 67–72.
- 593 Bepler, T.; Berger, B. Learning the protein language: Evolution, structure, and function.
594 *Cell systems* **2021**, *12*, 654–669.
- 595 Lin, Z.; Akin, H.; Rao, R.; Hie, B.; Zhu, Z.; Lu, W.; Smetanin, N.; Verkuil, R.; Kabeli, O.;
596 Shmueli, Y.; Dos Santos Costa, A.; Fazel-Zarandi, M.; Sercu, T.; Candido, S.; Rives, A.
597 Evolutionary-scale prediction of atomic-level protein structure with a language model.
598 *Science* **2023**, *379*, 1123–1130.
- 599 Elnaggar, A.; Heinzinger, M.; Dallago, C.; Rehawi, G.; Wang, Y.; Jones, L.; Gibbs, T.;
600 Feher, T.; Angerer, C.; Steinegger, M.; Bhowmik, D.; Rost, B. ProtTrans: Toward
601 Understanding the Language of Life Through Self-Supervised Learning. *IEEE Trans-
602 actions on Pattern Analysis and Machine Intelligence* **2022**, *44*, 7112–7127.
- 603 Kalakoti, Y.; Wallner, B. AFsample2: Predicting multiple conformations and ensembles
604 with AlphaFold2. *bioRxiv* **2024**, 2024–05.
- 605 Wayment-Steele, H. K.; Ojoawo, A.; Otten, R.; Apitz, J. M.; Pitsawong, W.;
606 Hömberger, M.; Ovchinnikov, S.; Colwell, L.; Kern, D. Predicting multiple confor-
607 mations via sequence clustering and AlphaFold2. *Nature* **2023**, 1–3.
- 608 Del Alamo, D.; Sala, D.; Mchaourab, H. S.; Meiler, J. Sampling alternative conformational
609 states of transporters and receptors with AlphaFold2. *Elife* **2022**, *11*, e75751.
- 610 Stein, R. A.; Mchaourab, H. S. SPEACH_AF: Sampling protein ensembles and con-
611 formational heterogeneity with AlphaFold2. *PLOS Computational Biology* **2022**, *18*,
612 e1010483.

- 613 Porter, L. L.; Artsimovitch, I.; Ramírez-Sarmiento, C. A. Metamorphic proteins and how
614 to find them. *Current opinion in structural biology* **2024**, *86*, 102807.
- 615 Chakravarty, D.; Porter, L. L. AlphaFold2 fails to predict protein fold switching. *Protein*
616 *Science* **2022**, *31*, e4353.
- 617 Chakravarty, D.; Schafer, J. W.; Chen, E. A.; Thole, J.; Porter, L. AlphaFold2 has more
618 to learn about protein energy landscapes. *bioRxiv* **2023**, 2023–12.
- 619 Hollingsworth, S. A.; Dror, R. O. Molecular dynamics simulation for all. *Neuron* **2018**,
620 *99*, 1129–1143.
- 621 Chen, H.; Roux, B.; Chipot, C. Discovering reaction pathways, slow variables, and com-
622 mitter probabilities with machine learning. *Journal of Chemical Theory and Compu-*
623 *tation* **2023**, *19*, 4414–4426.
- 624 Noé, F.; Tkatchenko, A.; Müller, K.-R.; Clementi, C. Machine learning for molecular
625 simulation. *Annual review of physical chemistry* **2020**, *71*, 361–390.
- 626 Belkacemi, Z.; Gkeka, P.; Lelièvre, T.; Stoltz, G. Chasing collective variables using au-
627 toencoders and biased trajectories. *Journal of chemical theory and computation* **2021**,
628 *18*, 59–78.
- 629 Bonati, L.; Piccini, G.; Parrinello, M. Deep learning the slow modes for rare events
630 sampling. *Proceedings of the National Academy of Sciences* **2021**, *118*, e2113533118.
- 631 Wang, Y.; Ribeiro, J. M. L.; Tiwary, P. Machine learning approaches for analyzing and
632 enhancing molecular dynamics simulations. *Current opinion in structural biology* **2020**,
633 *61*, 139–145.
- 634 Ribeiro, J. M. L.; Bravo, P.; Wang, Y.; Tiwary, P. Reweighted autoencoded variational
635 Bayes for enhanced sampling (RAVE). *The Journal of chemical physics* **2018**, *149*.
- 636 others,, et al. Predicting equilibrium distributions for molecular systems with deep learn-
637 ing. *Nature Machine Intelligence* **2024**, 1–10.
- 638 Lu, J.; Zhong, B.; Tang, J. Score-based enhanced sampling for protein molecular dy-
639 namics. ICML 2023 Workshop on Structured Probabilistic Inference {&} Generative
640 Modeling. 2023.
- 641 Ramaswamy, V. K.; Musson, S. C.; Willcocks, C. G.; Degiacomi, M. T. Deep learning pro-
642 tein conformational space with convolutions and latent interpolations. *Physical Review*
643 *X* **2021**, *11*, 011052.
- 644 Noé, F.; Olsson, S.; Köhler, J.; Wu, H. Boltzmann generators: Sampling equilibrium
645 states of many-body systems with deep learning. *Science* **2019**, *365*, eaaw1147.
- 646 others,, et al. Accurate structure prediction of biomolecular interactions with AlphaFold
647 3. *Nature* **2024**, 1–3.
- 648 Jing, B.; Erives, E.; Pao-Huang, P.; Corso, G.; Berger, B.; Jaakkola, T. Eigen-
649 Fold: Generative Protein Structure Prediction with Diffusion Models. *arXiv preprint*
650 *arXiv:2304.02198* **2023**,

- 651 Jing, B.; Berger, B.; Jaakkola, T. AlphaFold meets flow matching for generating protein
652 ensembles. *arXiv preprint arXiv:2402.04845* **2024**,
- 653 Hayward, S.; Go, N. Collective variable description of native protein dynamics. *Annual*
654 *review of physical chemistry* **1995**, *46*, 223–250.
- 655 Grudinin, S.; Laine, E.; Hoffmann, A. Predicting protein functional motions: an old recipe
656 with a new twist. *Biophysical journal* **2020**, *118*, 2513–2525.
- 657 Hoffmann, A.; Grudinin, S. NOLB: Nonlinear rigid block normal-mode analysis method.
658 *Journal of chemical theory and computation* **2017**, *13*, 2123–2134.
- 659 Laine, E.; Grudinin, S. HOPMA: Boosting protein functional dynamics with colored
660 contact maps. *The Journal of Physical Chemistry B* **2021**, *125*, 2577–2588.
- 661 Berman, H. M.; Westbrook, J.; Feng, Z.; Gilliland, G.; Bhat, T. N.; Weissig, H.;
662 Shindyalov, I. N.; Bourne, P. E. The protein data bank. *Nucleic acids research* **2000**,
663 *28*, 235–242.
- 664 Ramelot, T. A.; Tejero, R.; Montelione, G. T. Representing structures of the multiple
665 conformational states of proteins. *Current Opinion in Structural Biology* **2023**, *83*,
666 102703.
- 667 Bryant, P. Structure prediction of alternative protein conformations. *bioRxiv* **2023**, 2023–
668 09.
- 669 Trehwella, J. Recent advances in small-angle scattering and its expanding impact in
670 structural biology. *Structure* **2022**, *30*, 15–23.
- 671 Martel, A.; Gabel, F. Time-resolved small-angle neutron scattering (TR-SANS) for struc-
672 tural biology of dynamic systems: Principles, recent developments, and practical guide-
673 lines. *Methods in enzymology* **2022**, *677*, 263–290.
- 674 Flechsig, H.; Ando, T. Protein dynamics by the combination of high-speed AFM and
675 computational modeling. *Current Opinion in Structural Biology* **2023**, *80*, 102591.
- 676 Wankowicz, S.; Fraser, J. Comprehensive Encoding of Conformational and Compositional
677 Protein Structural Ensembles through mmCIF Data Structure. *ChemRxiv* **2023**,
- 678 Ellaway, J. I.; Anyango, S.; Nair, S.; Zaki, H. A.; Nadzirin, N.; Powell, H. R.; Gut-
679 manas, A.; Varadi, M.; Velankar, S. Identifying Protein Conformational States in the
680 PDB and Comparison to AlphaFold2 Predictions. *bioRxiv* **2023**, 2023–07.
- 681 Varadi, M. et al. PDBe and PDBe-KB: Providing high-quality, up-to-date and integrated
682 resources of macromolecular structures to support basic and applied research and ed-
683 ucation. *Protein Science* **2022**, *31*.
- 684 Modi, V.; Dunbrack Jr, R. L. Kincore: a web resource for structural classification of
685 protein kinases and their inhibitors. *Nucleic Acids Research* **2022**, *50*, D654–D664.
- 686 Parker, M. I.; Meyer, J. E.; Golemis, E. A.; Dunbrack Jr, R. L. Delineating the RAS
687 conformational landscape. *Cancer research* **2022**, *82*, 2485–2498.

- 688 Tordai, H.; Suhajda, E.; Sillitoe, I.; Nair, S.; Varadi, M.; Hegedus, T. Comprehensive
689 collection and prediction of ABC transmembrane protein structures in the AI era of
690 structural biology. *International Journal of Molecular Sciences* **2022**, *23*, 8877.
- 691 Pándy-Szekeres, G.; Caroli, J.; Mamyrbekov, A.; Kermani, A. A.; Keserű, G. M.; Koois-
692 tra, A. J.; Gloriam, D. E. GPCRdb in 2023: state-specific structure models using
693 AlphaFold2 and new ligand resources. *Nucleic Acids Research* **2023**, *51*, D395–D402.
- 694 Lombard, V.; Grudinin, S.; Laine, E. Explaining Conformational Diversity in Protein
695 Families through Molecular Motions. *Scientific Data* **2024**, *11*, 752.
- 696 others,, et al. Simulating 500 million years of evolution with a language model. *bioRxiv*
697 **2024**, 2024–07.
- 698 Heinzinger, M.; Weissenow, K.; Sanchez, J. G.; Henkel, A.; Steinegger, M.; Rost, B.
699 ProstT5: Bilingual language model for protein sequence and structure. *bioRxiv* **2023**,
700 2023–07.
- 701 Kim, D.-H.; Kang, S.-M.; Baek, S.-M.; Yoon, H.-J.; Jang, D. M.; Kim, H. S.; Lee, S. J.;
702 Lee, B.-J. Role of PemI in the Staphylococcus aureus PemIK toxin–antitoxin complex:
703 PemI controls PemK by acting as a PemK loop mimic. *Nucleic Acids Research* **2022**,
704 *50*, 2319–2333.
- 705 Machen, A. J.; Fisher, M. T.; Freudenthal, B. D. Anthrax toxin translocation complex
706 reveals insight into the lethal factor unfolding and refolding mechanism. *Scientific Re-*
707 *ports* **2021**, *11*, 13038.
- 708 Anderson, D. M.; Sheedlo, M. J.; Jensen, J. L.; Lacy, D. B. Structural insights into the
709 transition of Clostridioides difficile binary toxin from prepore to pore. *Nature microbi-*
710 *ology* **2020**, *5*, 102–107.
- 711 Dhanasingh, I.; Choi, E.; Lee, J.; Lee, S. H.; Hwang, J. Functional and structural char-
712 acterization of Deinococcus radiodurans R1 MazEF toxin-antitoxin system, Dr0416-
713 Dr0417. *Journal of Microbiology* **2021**, *59*, 186–201.
- 714 Chandravanshi, M.; Samanta, R.; Kanaujia, S. P. Conformational trapping of a β -
715 glucosides-binding protein unveils the selective two-step ligand-binding mechanism of
716 ABC importers. *Journal of molecular biology* **2020**, *432*, 5711–5734.
- 717 Hou, X.; Burstein, S. R.; Long, S. B. Structures reveal opening of the store-operated
718 calcium channel Orai. *Elife* **2018**, *7*, e36758.
- 719 Amadei, A.; Ceruso, M. A.; Di Nola, A. On the convergence of the conformational co-
720 ordinates basis set obtained by the essential dynamics analysis of proteins' molecular
721 dynamics simulations. *Proteins: Structure, Function, and Bioinformatics* **1999**, *36*,
722 419–424.
- 723 Leo-Macias, A.; Lopez-Romero, P.; Lupyan, D.; Zerbino, D.; Ortiz, A. R. An analysis of
724 core deformations in protein superfamilies. *Biophysical journal* **2005**, *88*, 1291–1299.
- 725 David, C. C.; Jacobs, D. J. Characterizing protein motions from structure. *Journal of*
726 *Molecular Graphics and Modelling* **2011**, *31*, 41–56.

- 727 others,, et al. A structural dendrogram of the actinobacteriophage major capsid proteins
728 provides important structural insights into the evolution of capsid stability. *Structure*
729 **2023**, *31*, 282–294.
- 730 Joosten, R. P.; Long, F.; Murshudov, G. N.; Perrakis, A. The PDB-REDO server for
731 macromolecular structure model optimization. *IUCrJ* **2014**, *1*, 213–220.
- 732 Devlin, J.; Chang, M.-W.; Lee, K.; Toutanova, K. Bert: Pre-training of deep bidirectional
733 transformers for language understanding. *arXiv preprint arXiv:1810.04805* **2018**,
- 734 Suzek, B. E.; Wang, Y.; Huang, H.; McGarvey, P. B.; Wu, C. H.; Consortium, U. UniRef
735 clusters: a comprehensive and scalable alternative for improving sequence similarity
736 searches. *Bioinformatics* **2015**, *31*, 926–932.
- 737 Suzek, B. E.; Huang, H.; McGarvey, P.; Mazumder, R.; Wu, C. H. UniRef: comprehensive
738 and non-redundant UniProt reference clusters. *Bioinformatics* **2007**, *23*, 1282–1288.
- 739 Consortium, T. U. UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic Acids*
740 *Research* **2022**, *51*, D523–D531.
- 741 Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.;
742 Polosukhin, I. Attention is all you need. *Advances in neural information processing*
743 *systems* **2017**, *30*.
- 744 others,, et al. MGnify: the microbiome sequence data analysis resource in 2023. *Nucleic*
745 *Acids Research* **2023**, *51*, D753–D759.
- 746 Olsen, T. H.; Boyles, F.; Deane, C. M. Observed Antibody Space: A diverse database of
747 cleaned, annotated, and translated unpaired and paired antibody sequences. *Protein*
748 *Science* **2022**, *31*, 141–146.
- 749 Barrio-Hernandez, I.; Yeo, J.; Jänes, J.; Mirdita, M.; Gilchrist, C. L.; Wein, T.;
750 Varadi, M.; Velankar, S.; Beltrao, P.; Steinegger, M. Clustering predicted structures at
751 the scale of the known protein universe. *Nature* **2023**, *622*, 637–645.
- 752 Varadi, M. et al. AlphaFold Protein Structure Database: massively expanding the struc-
753 tural coverage of protein-sequence space with high-accuracy models. *Nucleic Acids Re-*
754 *search* **2021**, *50*, D439–D444.
- 755 Van Kempen, M.; Kim, S. S.; Tumescheit, C.; Mirdita, M.; Lee, J.; Gilchrist, C. L.;
756 Söding, J.; Steinegger, M. Fast and accurate protein structure search with Foldseek.
757 *Nature Biotechnology* **2024**, *42*, 243–246.
- 758 Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.;
759 Liu, P. J. Exploring the limits of transfer learning with a unified text-to-text trans-
760 former. *Journal of machine learning research* **2020**, *21*, 1–67.
- 761 Steinegger, M.; Mirdita, M.; Söding, J. Protein-level assembly increases protein sequence
762 recovery from metagenomic samples manifold. *Nature methods* **2019**, *16*, 603–606.
- 763 LeCun, Y.; Bengio, Y.; Hinton, G. Deep learning. *nature* **2015**, *521*, 436–444.
- 764 others,, et al. Rectifier nonlinearities improve neural network acoustic models. Proc. icml.
765 2013; p 3.

- 766 Srivastava, N.; Hinton, G.; Krizhevsky, A.; Sutskever, I.; Salakhutdinov, R. Dropout: a
767 simple way to prevent neural networks from overfitting. *The journal of machine learning*
768 *research* **2014**, *15*, 1929–1958.
- 769 others,, et al. Pytorch: An imperative style, high-performance deep learning library. *Ad-*
770 *vances in neural information processing systems* **2019**, *32*.
- 771 Gower, J. C.; Dijkstrahuis, G. B. *Procrustes problems*; OUP Oxford, 2004; Vol. 30.
- 772 Schönemann, P. H. A generalized solution of the orthogonal procrustes problem. *Psy-*
773 *chometrika* **1966**, *31*, 1–10.
- 774 Kingma, D. P.; Ba, J. Adam: A method for stochastic optimization. *arXiv preprint*
775 *arXiv:1412.6980* **2014**,
- 776 Brüsweiler, R. Collective protein dynamics and nuclear spin relaxation. *The Journal of*
777 *Chemical Physics* **1995**, *102*, 3396–3403.
- 778 Tama, F.; Sanejouand, Y. H. Conformational change of proteins arising from normal
779 mode calculations. *Protein Engineering* **2001**, *14*, 1–6.