

# 1 Scalable design of repeat protein structural dynamics via 2 probabilistic coarse-grained models

3 Seeralan Sarvaharman<sup>1</sup>, Timon E. Neary<sup>2</sup>, Thomas E. Gorochowski<sup>1,3,\*</sup>, and Fabio  
4 Parmeggiani<sup>2,3,4,\*</sup>

5 <sup>1</sup> School of Biological Sciences, University of Bristol, 24 Tyndall Avenue, Bristol, BS8 1TQ, UK

6 <sup>2</sup> School of Biochemistry, University of Bristol, University Walk, Bristol, BS8 1TD, UK

7 <sup>3</sup> BrisEngBio, School of Chemistry, University of Bristol, Cantock's Close, Bristol, BS8 1TS, UK

8 <sup>4</sup> School of Chemistry, University of Bristol, Cantock's Close, Bristol, BS8 1TS, UK

9 \* These authors should be considered as joint senior author with equal contribution

10 Correspondence should be addressed to S.S. ([s.sarvaharman@bristol.ac.uk](mailto:s.sarvaharman@bristol.ac.uk)), T.E.G.

11 ([thomas.gorochowski@bristol.ac.uk](mailto:thomas.gorochowski@bristol.ac.uk)) and F.P. ([fabio.parmeggiani@bristol.ac.uk](mailto:fabio.parmeggiani@bristol.ac.uk))

## 12 ABSTRACT

13 Computational protein design has emerged as a powerful tool for creating proteins with novel functionalities.  
14 However, most existing methods ignore structural dynamics even though they are known to play a central  
15 role in many protein functions. Furthermore, methods like molecular dynamics that are able to simulate  
16 protein movements are computationally demanding and do not scale for the design of even moderately sized  
17 proteins. Here, we develop a probabilistic coarse-grained model to overcome these limitations and support  
18 the design of the structural dynamics of modular repeat proteins. Our model allows us to rapidly calculate the  
19 probability distribution of structural conformations of large modular proteins, enabling efficient screening  
20 of design candidates based on features of their dynamics. We demonstrate this capability by exploring  
21 the design landscape of 4–6 module repeat proteins. We assess the flexibility, curvature and multi-state  
22 potential of over 65,000 protein variants and identify the roles that particular modules play in controlling  
23 these features. Although our focus here is on protein design, the methods developed are easily generalised to  
24 any modular structure (e.g., DNA origami), offering a means to incorporate dynamics into diverse biological  
25 design workflows.

## 26 INTRODUCTION

27 The structural dynamics of proteins play a crucial role in their function and contribution to a wide variety  
28 of biomolecular processes<sup>1</sup>. Examples range from the active transport of molecules<sup>2</sup>, to the sensing of  
29 stimuli<sup>3</sup>. In the field of synthetic biology, computational protein design has emerged as a powerful tool  
30 for creating new proteins with desired functionalities. It has been used to support the design of florescence

31 activating proteins<sup>4</sup>, triose-phosphate isomerase (TIM) barrels<sup>5</sup>, proteins that can triggering immune re-  
32 sponses<sup>6</sup>, enzymes for catalysis<sup>7</sup> and protein switches<sup>8</sup>, to name but a few. However, while the central  
33 role of protein structural dynamics is well known<sup>9</sup>, when it comes to engineering *de novo* proteins, their  
34 dynamics for the most part have been neglected.

35 The most detailed predictions of protein dynamics are generated using molecular dynamics (MD) sim-  
36 ulations<sup>10-12</sup>. These provide atomistic detail and can capture the complex motions that proteins exhibit.  
37 Unfortunately, MD simulations are still time-consuming to perform and require substantial computational  
38 resources to be performed at scale, especially when simulating large proteins. Moreover, to achieve accurate  
39 prediction of dynamics, extensive sampling and long simulation times are necessary, limiting their applica-  
40 tion for selecting candidate designs from large variant libraries or to support rapid iterative design cycles.  
41 Some of these difficulties have been partially addressed by exploiting hybrid modelling approaches<sup>13,14</sup>,  
42 however, scalability issues remain.

43 To overcome the computational limitations of molecular simulations, coarse-grained models have been  
44 developed<sup>15-17</sup>. These use simplified representations of proteins, typically abstracting multiple atoms as  
45 a single interaction site. This reduces the degrees of freedom in the model and enables simulations over  
46 much longer timescales. Coarse-grained models in some cases have been found to capture the essential  
47 structural and dynamical features needed for design tasks, while significantly reducing the computational  
48 demands<sup>18,19</sup>.

49 A common coarse-grained approach for the prediction of protein dynamics is the elastic network model  
50 (ENM)<sup>20,21</sup>. ENMs approximate a protein structure as a network of interconnected springs, where each  
51 spring represents an interaction between two residues<sup>22</sup>. The model captures the collective motions of  
52 the protein by considering the harmonic vibrations around the equilibrium positions<sup>23</sup>. ENMs provide a  
53 simplified representation of protein dynamics and have been successful in modelling global, low-frequency  
54 motions and functionally important motions away from thermal fluctuations. However, they do not always  
55 fully capture higher frequency movements or the details of local interactions that result from non-collective  
56 motions<sup>24</sup>. Furthermore, the number of spring constants that must initially be fit changes depending on the  
57 number of amino acids present. This makes it difficult to rapidly evaluate similar candidates that differ in  
58 length, or whose model parameters are drastically different. Therefore, ENM models are typically unsuitable  
59 for iterative design workflows, where many diverse designs needs to be quickly evaluated at each cycle.

60 In this work, we aim to overcome these limitations and strike a balance between the computational effi-  
61 ciency of a coarse-grained representation and the ability for it to capture key protein dynamics. We focus on  
62 the simulation of tandem repeat protein domains, as they can reach hundreds of amino acids in size, resulting  
63 in large global motions as well as confined local motions<sup>25</sup>. Repeat proteins are found widely in nature<sup>26,27</sup>

64 and perform diverse biological functions<sup>28–30</sup>. They are characterised by the presence of repetitive struc-  
65 tural motifs/modules, which offers several advantages for modelling and *de novo* protein design. First, their  
66 modular nature greatly simplifies the prediction of tertiary structures, enabling the more rational design of  
67 novel sequences based on known repeat modules<sup>31</sup> and the combination of different repeats<sup>32</sup>. Second, and  
68 most importantly, repeat proteins exhibit more predictable structural dynamics within each module<sup>33–35</sup>,  
69 making them a useful platform upon which dynamics can be effectively modelled and exploited in the de-  
70 sign process<sup>36</sup>. By capitalising on the repetitive architecture of these proteins, we are able to construct a  
71 coarse-grained model that can efficiently propagate expected structural dynamics through the chain of mod-  
72 ules making up a protein (**Figure 1**). We demonstrate how the speed of our model allows us to predict, score  
73 and extract profiles of protein dynamics in a modular design landscape, offering a means to quickly and re-  
74 liably discover candidates with desired structural and dynamical characteristics. In addition, we show how  
75 this comprehensive view of a design space can provide valuable information regarding the flexibility and  
76 responsiveness of candidate modules to the sequence context, and offer insights into how specific modules  
77 are likely to affect overall features of a larger repeat protein. As protein design moves towards applications  
78 that require the careful crafting of conformational changes in protein structure, our model provides a means  
79 to assess such features and support engineering workflows that place dynamics at the forefront.

## 80 **RESULTS**

### 81 **Coarse-grained model of repeat protein dynamics**

82 The input to our model is a repeat protein library<sup>37</sup> consisting of a set of protein modules (with each mod-  
83 ule comprising at least two repeats) and a connectivity matrix that defines the rules for assembling larger  
84 constructs, i.e., for a given module, which other modules are compatible and can be directly connected  
85 (**Figure 1a**, step 1). Note that these rules do not necessarily commute, such that “module B can follow mod-  
86 ule A” does not imply “module A can follow module B”. This library defines the overall accessible design  
87 space. However, it can be extended at any time by adding further modules and connectivity rules. We chose  
88 to use an existing repeat protein library that contains 34 modules that are on average 180 Å long, covering a  
89 wide range of structures<sup>38,39</sup>.

90 In addition to the repeat protein library, we also require a dynamics database that can be used by the  
91 model to predict the movement of larger multi-module proteins (**Figure 1a**, step 2). The database consists  
92 of a large number of conformational snapshots for all possible combinations of three-module constructs.  
93 For our library, this equates to a total of 644 unique constructs that were 320 to 794 amino acids long, and  
94 100 conformational snapshots for each. These snapshots can be identified with minima in the rough energy  
95 landscape of the protein, and the movement of the protein is driven by thermal fluctuations, which on a slow

96 enough timescale cause jumps from one minima to another. These snapshots therefore allow us to infer the  
97 steady-state dynamics of proteins using a probability distribution of atomistic positions.

98 Obtaining this information via experimental methods such as hydrogen-deuterium exchange via Nuclear  
99 Magnetic Resonance (NMR) or mass spectrometry (HDX-MS), is infeasible due to the total number of  
100 unique constructs in our library and the size of the molecules. Using computational methods like molecular  
101 dynamics (MD), simulations were also not feasible due to the size of the proteins. We therefore chose  
102 to use conformational snapshots generated using the relax protocol from the Rosetta modelling suite<sup>40</sup>  
103 (**Methods**). We based our database on three module constructs due to this being the smallest number of  
104 modules where contextual effect of different neighbours on a given module can be observed. While it  
105 is possible to utilise higher-order contextual effects by using constructs with more than three modules,  
106 obtaining snapshots becomes computationally challenging due to the larger number of constructs that must  
107 be assessed and the increased number of residues per construct.

108 For the model to efficiently use the information within the dynamics database, we generated coarse-  
109 grained descriptions that simplify the propagation of dynamical information throughout larger multi-module  
110 proteins (**Figure 1a**, step 3). We chose key anchor points, specifically the ends of  $\alpha$ -helices, to capture the  
111 orientation of the module and then define the mean of these key anchor points as the centroid of the module.  
112 The locations of the centroids and the anchor points change depending on the conformation. Therefore, we  
113 use the conformational landscape in the dynamics database to build a probability density estimate for the  
114 locations of the centroids and their anchor points. As these descriptions are analogous to position vectors  
115 when the dynamics are neglected, we refer to them as probabilistic vectors. For each unique triplet of  
116 modules ( $a \cdot b \cdot c$ ), the coarse-grained model is then generated as follows. We first perform a rigid body  
117 transformation on all of the conformations such that the centroid of the central module  $b$  ( $\mathbf{c}_b$ ) is located at  
118 the origin ( $\mathbf{0}$ ), and rotated appropriately such that any rotational symmetry from the placement of modules  
119 is removed. Using the conformational data, we then calculate the probability distribution  $P(\mathbf{c}_a | \mathbf{c}_b = \mathbf{0})$   
120 and  $P(\mathbf{c}_c | \mathbf{c}_b = \mathbf{0})$ , which captures the steady-state occupation probability density of module  $a$  relative to  
121 module  $b$ , and  $c$  relative to module  $b$ , respectively. For module  $b$ , we can also track an arbitrary number of  
122 reference points  $\mathbf{r}_i$  (e.g., start and ends of alpha helices) via  $P(\mathbf{r}_i | \mathbf{c}_b = \mathbf{0})$ , which describes the steady-state  
123 occupation probability density of the  $i^{\text{th}}$  reference point in module  $b$  relative to the centroid of module  $b$ .  
124 These probability density functions are estimated by fitting a Gaussian mixture that can be stored efficiently  
125 using the parameters of the mixture (i.e., means, covariances and weights of the constituent components).  
126 We precompute these parameters for all of the module triplets and use them for efficiently estimating the  
127 dynamics of larger constructs.

128 Finally, to predict the dynamics of arbitrarily sized modular proteins, we use this description to estimate

129 the steady-state occupation probability density of the  $k^{\text{th}}$  centroid in a many-module construct by identifying  
130 the constituting triplets that make up the larger construct, and perform the convolutions

$$P(\mathbf{c}_k|\mathbf{c}_1 = \mathbf{0}) = \int d\mathbf{c}_{k-1} \cdots \int d\mathbf{c}_2 P(\mathbf{c}_k|\mathbf{c}_{k-1}) \cdots P(\mathbf{c}_2|\mathbf{c}_1 = \mathbf{0}) \quad (1)$$

131 when  $k \geq 3$  and where we condition the location of the centroid of the first module to be at the origin  
132 (**Figure 1a**, step 4). To illustrate the procedure, we can consider a three module construct, labelled 1, 2 and  
133 3. With the centroid of the first module  $\mathbf{c}_1$  centred at the origin,  $P(\mathbf{c}_2|\mathbf{c}_1 = \mathbf{0})$  describes the position of the  
134 second centroid  $\mathbf{c}_2$  relative to the first. Similarly,  $P(\mathbf{c}_3|\mathbf{c}_2 = \mathbf{0})$  describes the position of the third centroid  
135  $\mathbf{c}_3$  relative to a fixed  $\mathbf{c}_2$  at the origin. In order to obtain the distribution of the position of  $\mathbf{c}_3$  relative to the  
136 the fixed  $\mathbf{c}_1$  at the origin, one needs to perform a convolution of the distributions. The probability density  
137 estimate for the reference points in the  $k^{\text{th}}$  module can also then be computed via:

$$P(\mathbf{r}_i|\mathbf{c}_1 = \mathbf{0}) = \int d\mathbf{c}_k P(\mathbf{r}_i|\mathbf{c}_k)P(\mathbf{c}_k|\mathbf{c}_1 = \mathbf{0}). \quad (2)$$

138 Together, these equations allow us to propagate the local movements captured in our dynamics database  
139 through to larger constructs and estimate diverse structural dynamics from the centroids and reference points.

140 This algorithm was implemented in a package called Dynamo (**Data Availability**). Dynamo is a native  
141 Python library written in Rust to ensure reliable and high-performance model generation and simulation.  
142 It also includes additional helper functions to simplify the creation of the dynamics database, the ability to  
143 define complex multi-module constructs beyond chains (e.g., star-like proteins), visualisation tools to better  
144 understand the inferred protein dynamics, and the ability to export data in standard Protein Data Bank (PDB)  
145 format for use in other tools.

## 146 **Model validation**

147 To verify the accuracy of our model, we assessed the differences between the dynamics predicted from our  
148 coarse-grained model with those extracted from the conformations obtained from the Rosetta relax protocol.  
149 We chose to consider a diverse set of 14 modular proteins where 10 were homogeneous containing 9 modules  
150 of the same type, while 4 were heterogeneous containing 4 modules of different types. These proteins ranged  
151 in size from 701 to 840 amino acids long.

152 To compare the qualitative agreement between our model and the Rosetta relax data, we computed  
153 the displacements,  $\mathbf{r}$ , between the specific centroids of the conformational samples and the mean centroid  
154 location across them all, and compared the distribution of its magnitude  $|\mathbf{r}|$ . This distribution captures both  
155 the overall magnitude of any movement, as well as its shape. For example, if the centroid was distributed  
156 on the surface of a sphere of radius  $R$ , then the distribution of  $|\mathbf{r}|$ , would reduce to the Dirac-delta function  
157  $P(|\mathbf{r}|) = \delta(|\mathbf{r}| - R)$ . As the point cloud of  $\mathbf{r}$  becomes more complex in shape, the distribution  $P(|\mathbf{r}|)$  will

158 exhibit more complex features.

159 For each modular protein construct, we plotted the distribution  $P(|r|)$ , for the 3<sup>rd</sup>, 6<sup>th</sup>, and 9<sup>th</sup> cen-  
160 troids for the homogeneous constructs, while the 2<sup>nd</sup>, 3<sup>rd</sup>, and 4<sup>th</sup> centroids for the heterogeneous constructs  
161 (**Figure 2a**). In most cases, we found model predictions agreed well with the Rosetta relax data (i.e., the  
162 modes of the distribution coincided with each other). The main exception was the D4 construct, where  
163 the model predicted more movement than expected. A potential reason for this disagreement could be the  
164 larger number of rare conformations of the D4-D4-D4 triplet in the dynamics database used to parameterise  
165 the model. This would result in a wider exploration of the conformational landscapes and cause the model  
166 to infer larger movements that get propagated through the entire protein. For the H4 construct, we also  
167 overestimated the dynamics. However, the differences are smaller than those of the D4 construct.

168 To compare the distributions further, we computed the percentage error in the means of the distributions  
169 for each of the centroids (**Figure 2b**, top panel). Again, we found that the largest differences between the  
170 model and Rosetta relax data was for the D4 construct. However, the average error across all constructs and  
171 centroids was only 0.24% which further provides evidence for the accuracy of the inferred movements of  
172 the proteins. To better assess the similarity in the shape of the distributions, we also calculated the Earth  
173 mover's distance for each centroid in every protein construct (**Figure 2b**, bottom panel). We found that  
174 the average Earth mover's distance across all centroids and constructs was only 0.95 Å, suggesting good  
175 quantitative agreement despite the coarseness of our the model.

## 176 **Visualisation of protein structural dynamics**

177 During model validation, it became clear that it was difficult to assess subtle differences in protein dynamics  
178 due to the need to observe both structural and movement features of the data simultaneously. To help  
179 overcome this, we developed a new visualisation technique that captures the major dynamical features of  
180 each module in the context of the core alpha-helices making up the protein (**Figure 3**). The visualisation is  
181 created by first extracting the steady-state distributions of each centroid within a construct and calculating  
182 the covariances in their movement. For each module covariance, the direction and magnitude of the principal  
183 variances can be described as three mutually orthogonal vectors. These can then be used to generate 'fins',  
184 spanning either side of the mean centroid position and parallel to the principal variances. To aid comparison,  
185 we colour the fin at a particular point in relation to the total variance in the distribution of centroid locations,  
186 which corresponds to the magnitude of any movement. We also overlay each helix as a semi-transparent grey  
187 cylinder to further portray the protein's underlying structure. Using this visualisation technique, we could  
188 place the errors in the context of the protein size and shape and clearly see that our coarse-grained model  
189 was able to accurately capture the specific magnitude and direction of protein dynamics when compared to

190 the Rosetta relax data (**Figure 3**).

### 191 **Mapping out the design landscape of a repeat protein library**

192 The efficiency of our model enables us to predict the structural dynamics of large modular protein design  
193 spaces. To demonstrate this, we considered all possible 4-module protein designs using our modular protein  
194 database. We generated the structural dynamics of each of the 2,978 designs within this space and quantified  
195 two key features that covered the structural geometry and dynamics of each design. The first was the protein  
196 curvature  $\alpha$ , defined as the ratio between first and last centroid distance and corresponding arc length. The  
197 second was the flexibility  $\beta$ , which captures the overall extent of the dynamics (see details in methods).  
198 Plotting the density of designs in relation to these features provided us with a uniquely complete picture of  
199 the available characteristics of all 4-module proteins (**Figure 4a**).

200 Within the design space, we selected four examples with different properties and verified their qualitative  
201 agreement with data obtained from the Rosetta relax protocol. Our model predicted design I to be the most  
202 flexible (highest  $\beta$ ) and design IV the least (lowest  $\beta$ ). This was verified by the mean Root Mean Squared  
203 Deviation (RMSD) of the final centroid being the highest and lowest, respectively, and by distributions with  
204 the largest and narrowest distributions (**Figure 4a**, inset). Moreover, the model also captured the fact that  
205 designs II and III have very similar flexibility profiles corresponding to similar values of  $\beta$  and  $C_4$  RMSD  
206 distributions (**Figure 4a**, inset), whilst having very different sinuosity (**Figure 4b**).

### 207 **Proteins with multi-state potential**

208 The ability for a protein to adopt several distinct conformational states plays a crucial role in a wide variety  
209 of cellular processes spanning, the function of molecular motors<sup>41</sup> to improved catalysis<sup>42</sup>. From a design  
210 perspective, being able to suggest candidate protein designs with an inherent propensity for multi-state  
211 conformations would be valuable as a starting point for switchable and dynamic protein functions. While  
212 both globular and non-globular proteins can exhibit multi-state dynamics, repeat proteins represent a highly  
213 predictable and versatile platform to design multi-state systems. Due to the localised interactions, multi-state  
214 dynamics exhibited by a repeat protein can emerge from the behaviour of the individual modules.

215 We capitalised on this feature and developed a method to score candidate protein designs based on their  
216 potential for multi-state dynamics. Because multi-state dynamics are manifested as multi-modal probability  
217 distributions in the centroid positions, we could quantify the multi-modality of a construct by considering  
218 the probability density of the centroid of the final module. To do this, we first sampled from its probability  
219 density estimate. These samples were then split into training and testing sets using a  $k$ -fold validation  
220 scheme ( $k = 5$ ). We used the training set to fit a Gaussian mixture model with the number of components,  
221  $n$ , ranging from 1 to 5. For each fit, we obtained a score given by the likelihood of the appropriate test

222 sets, which was then normalised appropriately to construct a probability mass function over the number of  
223 components,  $p(n)$ . The mode of this distribution corresponded to the modes of the centroid density,  $\gamma_1$ . To  
224 score this aspect (i.e. presence of distinct multiple-modes), we then computed the reciprocal of the entropy  
225 of the probability mass function:

$$\gamma_2 = -\frac{1}{\sum_n p(n) \log p(n)}. \quad (3)$$

226 Low values of  $\gamma_2$  indicate a low confidence in the number of modes of the probability distribution of centroid,  
227 whereas high values a high confidence. Using this scoring we could rank candidates both on the number of  
228 modes present in their distribution and our confidence in this property.

229 To demonstrate this approach, we again considered all 4-module constructs and screened them for po-  
230 tential multi-state dynamics. Of our eight top scoring candidates, several showed bi-stable or multi-stable  
231 dynamics with H10 and H11 being the most promising (**Figure 5**). Many displayed very wide distributions  
232 with no clearly separated modes. While these candidates may not exhibit strong multi-state behaviour, the  
233 large range of movement (60 Å for H5) would provide an ideal starting point for establishing multi-stable  
234 dynamics using external stimuli, such as an additional peptide chain<sup>43</sup>.

### 235 **The effect of specific modules on construct rigidity**

236 The ability to generate the structural dynamics for entire modular protein design spaces, offers the ability to  
237 unravel the potential roles that individual modules might play more broadly across many different designs.  
238 A key feature that often needs to be controlled when designing *de novo* proteins is the rigidity of the final de-  
239 sign. This can be quantified by calculating  $\beta$ , which captures the likely overall movement of a protein from  
240 the probabilistic conformational data (**Methods**). Given that the overall rigidity of a construct is determined  
241 by the modules present (e.g., some modules might stabilise the module while others might introduce more  
242 flexibility), it is possible to uncover the influence each module has when they are part of a larger construct.  
243 To do this, for a given module in an  $n$  module construct, all possible constructs can be separated into  $n + 1$   
244 subsets based on the frequency of that module within the construct (**Figure 6a**, left panel). For each of these  
245 subsets, a probability density function can be built based on the constructs present. To compare how the  
246 properties of these distributions change with the different counts of a particular module, we can then visu-  
247 alise the means and variances of the distributions in a two-dimensional parameter space plot (**Figure 6a**,  
248 right panel). The mean of the distribution (Mean  $\beta$ ) denotes the average flexibility of the population, while  
249 the variance of the distribution (Var  $\beta$ ) quantifies the heterogeneity. As the occurrences of a specific module  
250 increases, the trajectory in this parameter space determines the role that the module plays within the larger  
251 constructs. There are four key types of behaviour: (i) universally stabilising, (ii) universally destabilising,  
252 (iii) contextually stabilising, and (iv) contextually destabilising. Increasing the counts of a stabilising mod-



253 ule in a construct reduces its flexibility which causes the mean of the distribution  $P^n(\beta)$  to decrease. When  
254 this stabilisation effect is universal, adding more of the given module in a construct drives down the dy-  
255 namics regardless of the local context where the module is being added, or the presence of other modules  
256 which in turn causes the variance of  $\beta$  to also decrease. Conversely, when the stabilising is contextual, the  
257 local context and combination of other modules present play more of an important role in the dynamics of  
258 the construct. In other words, the flexibility sub-population remains highly heterogeneous with increases in  
259 module counts which yields less reduction in the variance of  $\beta$  when compared with the analogue of the  
260 whole population. On the other hand, increasing the counts of a destabilising module causes the flexibility  
261 of the constructs to also increase. When the effect is independent of other contextual factors then it is called  
262 strongly destabilising yielding decreases in the variance of  $\beta$ . Whereas the weakly destabilising are those  
263 that increase the flexibility of the construct, but only within specific contexts. It should be noted that while  
264 we have here focused on flexibility/rigidity of the constructs, this types of approach can be used for any  
265 feature that can be calculated from the structural dynamics.

266 To quantify the role that modules had in relation to protein rigidity, we generated the structural dynamics  
267 for all possible 6-module constructs and assessed the parameter space plots for 18 modules in our database  
268 (**Figure 6b**). These showed a broad range of behaviours across the modules. Most prominent was a universal  
269 stabilising effect, displayed by many of the modules and most prominently by D49, D54, D18 and D14.  
270 D79 showed a less universal stabilising effect, while several modules displayed non-uniform behaviours  
271 (e.g., D4, D53 and D64). We also found that 4 of the modules were universally destabilising (D14.j1.D54,  
272 D14.j1.D79, D14.j2.D54 and D14.j4.D79), with D14.j4.D79 having the strongest effect. Interestingly,  
273 while D14 is strongly stabilising with some contextual dependence when used as part of a junction module,  
274 D14.j2.D14, it becomes more universal with a weaker stabilising effect. Whereas when used as part of  
275 D14.j1.D14, it has no stabilising effect at all. The latter is due to the presence of a junction domain which  
276 hinders the packing of helices into tight conformations.

277 To explore this unusual feature of the D14 module further, we considered constructs consisting of  $k$   
278 consecutive repeats of D14 modules followed by  $6 - k$  consecutive repeats of D14.j1.D14 or D14.j2.D14.  
279 For each of these constructs, we compared the distributions of pairwise distances between all carbon alphas  
280 and normalised these to the largest distance in the smallest construct (i.e., six repeats of D14). We found that  
281 as the number of D14s decrease and the number of D14.j1.D14s increases, the distributions flatten out and a  
282 second prominent mode emerges (**Figure 6c**). This coincides with a less tight packing that ultimately leads  
283 to a less stable behaviour. In contrast, as the number of D14s decrease and the number of D14.j2.D14s  
284 increases, the mode increases but the variance of the distribution decreases, leading to more consistent  
285 packing and stable behaviour.

## 286 **Modelling multi-chain constructs**

287 So far, we have only considered single-chain repeats. However, building more complex structures quickly  
288 becomes infeasible, as many structures are not reducible to a single-chain design. Moreover, due to the  
289 repetitive nature of the sequences, long repetitive DNA molecules can be difficult to synthesise and large  
290 repeat proteins can be expressed in low yield. One approach that is observed in nature, and commonly used  
291 by protein engineers to overcome this limitation is to employ multiple chains that physically interact to form  
292 a larger structure.

293 Tree-like structures bring together two or more single-chain repeat domains at branching points, which  
294 act like ‘hubs’ within the structure. To facilitate the design of tree-like structures using Dynamo, we ex-  
295 tended its capabilities to allow for hub modules within the parts database. Each hub is able to connect  
296 together multiple chains of modules together, allowing them to act as branches in the overall tree structure.

297 To accommodate tree-like structures in our model, we exploited the fact that when two modules are  
298 coupled, they physically interact and are stuck together. In other words, they do not necessarily have to be  
299 part of a single chain. This allowed us to loosen our definition of a module to merely regions of a protein that  
300 can be reused in different designs. In reality, hubs could either be an expressible protein, or could emerge  
301 from the cross interactions between two or more linear chains of modules. In this second case, we offer  
302 the ability to subjectively define a region around the interaction as a module. In terms of the abstraction in  
303 the model, we make no distinction between either case and treat both similarly (i.e., we represent hubs in  
304 the same way as modules, but with that added ability of being able to connect to more than two modules).  
305 With this extension, we can represent an arbitrary tree-like structure using a set of modules and a list of  
306 connections for each.

307 In order to predict the dynamics of tree-like structures, we exploited the fact that once represented  
308 using our model, the tree-like structure naturally defines a Bayesian network. We can therefore arbitrarily  
309 select an anchor module relative to which movement of other modules will be predicted, and then propagate  
310 movements using Eq. (1) for all branches emanating from the selected anchor module. The additional  
311 assumption underlying this approach is that the branches downstream of a hub do not physically interact  
312 with each other, which is valid when the ends of the branches are sufficiently separated.

313 To test this functionality, we designed a star-like multi-chain protein where a D4\_C4\_G1 hub module is  
314 connected to four independent chains of four D4 modules (**Figure 7a,b**). We then predicted the structural  
315 dynamics for two different anchoring points at the central hub and at the end of one of the arms. As expected,  
316 this showed that anchoring at the central hub reduced the overall movements that could be achieved by the  
317 arms, while anchoring a single arm allowed for larger arm movements (**Figure 7c**).

## 318 DISCUSSION

319 In this study, we developed a coarse-grained modelling approach to facilitate dynamics-driven repeat protein  
320 design. Our method successfully captured the essential features of modular protein dynamics and allowed  
321 for the exploration of their conformational space in a computationally efficient manner. For moderately  
322 sized proteins (4 or 6 modules long), the ability to calculate the conformational probability distributions  
323 and associated analyses in milliseconds on a standard desktop computer allowed us to exhaustively explore  
324 the structural dynamics for all possible designs, covering over 65,000 variants. The ability to provide such  
325 extensive coverage in protein design space enabled us to better understand how our design space covers  
326 particular features of interest, e.g., curvature and movement (**Figure 4**), and unravel the role that individ-  
327 ual modules play in supporting the flexibility/rigidity of a resultant protein (**Figure 6**). Furthermore, the  
328 generality of our approach is not limited to single-chain repeat proteins. We show that a simple extension  
329 enables the prediction of structural dynamics of multi-chain, tree-like modular proteins (**Figure 7**) and the  
330 underlying mathematical model can accommodate any modular component for which samples of structural  
331 conformations can be gathered. The current work was focused on a library of compatible alpha helical mod-  
332 ules, purely because of their availability, but, as more modular designs become available, our method can  
333 be apply to any modular system, even beyond proteins. Similarly, any dataset capturing population dynam-  
334 ics, either experimentally obtained or generated through simulations, could be used for the description and  
335 analysis of modular systems.

336 A major advantage of our approach is that it allows us to capture the inherent flexibility of repeat pro-  
337 teins. Modular repeat proteins often exhibit structural plasticity, allowing them to adapt and interact with  
338 different ligands or partners. By evaluating the dynamics of repeat proteins using our coarse-grained model,  
339 we are able to observe conformational changes and fluctuations in protein structure. These insights provide  
340 valuable information about the flexibility of different regions within the repeat protein and how they might  
341 contribute to its function. Understanding the flexibility of repeat proteins is also crucial for designing pro-  
342 teins with adjustable properties or for engineering proteins that can undergo conformational changes (e.g.,  
343 upon binding to specific targets).

344 Our approach also revealed the multi-stability of many repeat proteins (**Figure 5**), which is a desirable  
345 property in many applications. Multi-stability refers to the ability of a protein to adopt multiple stable  
346 conformations or functional states. By exploring the conformational space of repeat proteins, we identified  
347 distinct energy minima corresponding to different conformations. This finding suggests that repeat proteins  
348 can exist in alternative stable states, potentially enabling them to switch between different functional states  
349 or adopt different binding configurations. Exploiting the multi-stability of repeat proteins opens up new

350 opportunities for designing protein-based switches, sensors, or molecular machines with programmable  
351 functionalities. Our feature extraction method for identifying multi-stable features could be further refined  
352 to allow for greater specificity.

353 With the advent of AlphaFold<sup>44</sup>, machine learning (ML) and generative artificial intelligence (AI) have  
354 become commonplace in protein design workflows<sup>45,46</sup>. While these approaches offer unprecedented ac-  
355 curacy in the prediction of protein structure from sequence alone, their use for the prediction of protein  
356 dynamics has been limited<sup>47</sup>. This stems in part from difficulties in generating the large training sets re-  
357 quired, although there have been recent efforts to overcome these issues<sup>48,49</sup>. A further challenge that  
358 remains is the high-computational cost of running ML models after training. While acceptable for small  
359 design spaces containing hundreds of possible designs, larger design spaces remain inaccessible due to the  
360 computational demands. An interesting future direction would be use ML to generate the conformational  
361 snapshots needed to parameterise the modules of the Dynamo model<sup>49</sup>. This would then offer the means to  
362 blend ML predictions at the protein module level, with Dynamo's efficiency in combining that data at the  
363 level of large single-chain repeat proteins or multi-protein assemblies.

364 The past decade has seen the design of *de novo* protein structures explode. Looking forward, the next  
365 frontier will be the design of protein dynamics and the push towards implementing complex molecular func-  
366 tions that require carefully choreographed structural changes over time. Tools like Dynamo will be crucial  
367 for accelerating our ability to practically explore the dynamics of repeat proteins and modular biological  
368 structures, supporting steps towards this goal.

## 369 **METHODS**

### 370 **Repeat protein library**

371 We use a reduced subset of an existing repeat protein library<sup>38,39</sup> consisting of 11 homo-modules, and 23  
372 junction modules that consists of two homo-modules interfaced together with a junction modules. The  
373 repeat protein library can be used as part of the Elfin<sup>37</sup> tool or directly from the data repository (i.e.,  
374 <https://github.com/Parmeggiani-Lab/elfin-data>) where atomistic position data for all  
375 of the modules are stored as PDB files in the compressed tarball `pdb_aligned.tar.bz2`.

### 376 **Coarse-grained representation**

377 Given that an atomic description is far too detailed for our purposes, it was important to find a simpler  
378 representation that is computationally tractable. An option would be to use the coordinates of the  $C_\alpha$  atoms,  
379 and it would be possible to move to this level of detail at some point in the future, but for this work we find an  
380 even coarser description of secondary structures works well. Specifically, we use the STRIDE algorithm<sup>50,51</sup>

381 as it exploits dihedral angle information in addition to the hydrogen bonds. While STRIDE can find all the  
382 structures in each of the conformations, the start and end locations can vary in the order of a few residues.  
383 Thus, to compare each of the conformations, we take the intersection of the start and end locations among  
384 all of the conformations. In other words, for a given alpha-helix identified, the residues that are part of it are  
385 given by all of the residues that are common across all of the conformations that have been attributed to the  
386 same alpha-helix.

### 387 **Generating the dynamics database**

388 To build the dynamics database we employed the Elfin software suite<sup>37</sup> to construct first all possible three  
389 module constructs. An exhaustive approach was used resulting in 644 repeat proteins. Each of these repeat  
390 proteins are single chains and they consist of three modules with capping repeats at both ends. Capping  
391 repeats are used to ensure solubility of the proteins when expressed, but in this work, they also prevented  
392 edge-dependent artifacts, such as opening of the terminal helices during relax. These proteins were further  
393 relaxed using the Cartesian relax protocol in the Rosetta relax application<sup>40,52,53</sup> to obtain a low energy ref-  
394 erence structure with packed side chains. Using these reference structures, an additional round of relaxation  
395 was used to obtain 100 conformations given as atomistic positions in PDB files. For this latter relaxation  
396 the FastRelax protocol was used and no conformations were rejected once they were obtained. All Rosetta  
397 related tasks were performed using Rosetta version 3.9 and ‘ref2015’ score function for the relaxations.

### 398 **Abstract representation and rules of combination**

399 Before we consider the dynamics of modules in a chain, we first construct a set of rules and axioms that are  
400 necessary to construct larger proteins. We start with the 2D representation before generalising to obtain the  
401 3D representation.

402 We have to design a unique local representation and connection rules that result in a unique “deter-  
403 ministic” solution. Let  $C = \{\mathcal{U}_1, \dots, \mathcal{U}_1\}$  be the set of modules and  $\oplus$  be a non-commutative operation  
404 to combine two modules. Suppose a module is defined in the following way, (a) a bounding box defined  
405 by a set of points that are measured relative to the centroid of mass (centroid). (b) a vector defining the  
406 location of the adjacent centroid w.r.t the current  $\mathbf{v}_n$ . (c) a vector connecting the current centre of mass to  
407 the previous,  $\mathbf{v}_p$ . With these properties, we define the operation  $\mathcal{M}_1 \oplus \mathcal{M}_2$  as the following. (1) Translate  
408  $\mathcal{M}_2$  such that its centroid,  $\mathbf{c}_2 = \mathbf{c}_1 + \mathbf{v}_{n:1 \rightarrow 2}$ . (2) Rotate,  $\mathcal{M}_2$  about its centroid such that  $\mathbf{v}_{p:1 \rightarrow 2}$  is parallel to  
409  $\mathbf{v}_{p:1 \rightarrow 2}$ . Notice that with just the operation (1), it is not sufficient to create a unique  $\mathcal{M}_3 = \mathcal{M}_1 \oplus \mathcal{M}_2$  as  $\mathcal{M}_2$   
410 can freely rotate about its centroid. Having the vector  $\mathbf{v}_p$  with which to align  $\mathbf{v}_n$  ensures that  $\mathbf{v}_n$  can only  
411 point in a single direction, i.e. removes radially non-uniqueness.

412 Generalising to 3D, we find the following problem. While,  $\mathcal{M}_2$  and  $\mathcal{M}_1$  are aligned with respect to the

413 centroid and the  $\mathbf{v}_n$ , we find that properties (a)-(c) and (1) and (2) are no longer sufficient to remove the  
414 problem of non-uniqueness. This is due to the fact that the module  $\mathcal{M}_2$  can freely rotate about the vector  
415  $\mathbf{v}_{n:1 \rightarrow 2}$ . The simplest way to amend this is by having another vector which is not parallel to the  $\mathbf{v}_n$

## 416 Defining units for database

417 For our particular protein database, we can distil all of these ideas from the previous section into a concise  
418 system. Note that for clarity, we use the term unit to define the mathematical representation of the protein,  
419 and reserve the term module to refer to the actual protein module.

420 From our database, each one of the 644 repeat proteins, defines a unique triplet of modules,  $({}_L\mathcal{M}_R)$ .  
421 Given that the position and movement of elements within a module depend on its context, i.e. neighbouring  
422 the modules  $L$  and  $R$ , it is necessary to define a module for each unique triplet of modules. In other words,  
423 our database of units has a size of 644, whereas the number of unique modules is only 34.

424 To define a unit, from a triplet  ${}_L C_R$  repeat protein, we employ the following steps. We first separate the  
425 helices that belong to  $L$ ,  $C$  and  $R$ . We then compute centroids,  $\mathbf{c}_L$ ,  $\mathbf{c}_C$ , and  $\mathbf{c}_R$  for the modules,  $L$ ,  $C$  and  $R$ ,  
426 respectively. We define a reference,  $\underline{\mathbf{R}}^{(L \rightarrow C)} = [\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3]$ . The vectors  $\mathbf{e}_i$  are normalised orthogonal vectors,  
427 with  $\mathbf{e}_1$  parallel to  $\mathbf{c}_C - \mathbf{c}_L$ , and  $\mathbf{e}_2$  parallel to  $\mathbf{e}_1 \times \mathbf{h}_L$ , where  $\mathbf{h}_C$  is the vector from the mean bottom of helices  
428 in module  $C$  to the mean top. Having define,  $\underline{\mathbf{R}}^{(L \rightarrow C)}$ , perform a rigid body transformation by multiplying  
429 all points in  $C$  and  $\mathcal{R}$  with the inverse of  $\underline{\mathbf{R}}^{(L \rightarrow C)}$ . We can then define another reference frame for  $\underline{\mathbf{R}}^{(C \rightarrow R)}$   
430 following the same constraints as before but modules  $C$  and  $R$ . To connect with the module  $R$ , we define  
431 the vector  $\mathbf{v}_n = \mathbf{c}_R - \mathbf{c}_C$ . Lastly, we can define any other reference points within the module  $C$ , e.g. those of  
432 a bounding convex hull. In summary, we compute  $\underline{\mathbf{R}}^{(C \rightarrow R)}$ ,  $\mathbf{v}_n$  and any bounding box points relative to the  
433 centroid of  $C$  and with the additional constraint that  $\underline{\mathbf{R}}^{(L \rightarrow C)} = \underline{\mathbf{I}}$ .

434 With this abstract representation, we define the non-commutative operation  $\oplus$ , for  $\mathcal{M}_1 \oplus \mathcal{M}_2$  as the fol-  
435 lowing. (1) We translate  $\mathcal{M}_2$  such that its centroid is at  $\mathbf{c}_1 + \mathbf{v}_{n:1 \rightarrow 2}$ . (2) Perform a rigid body transformation  
436 for all of the vectors and reference frame in  $\mathcal{M}_2$  by the  $\underline{\mathbf{R}}^{(C \rightarrow R)}$  reference frame of  $\mathcal{M}_1$ . With these two rules,  
437 we can construct any valid chain of modules.

## 438 Generalising the module abstraction to include dynamics

439 Having designed the module framework, we now move on to generalising the frame to capture the dynamics  
440 involved. Such a generalisation is relatively straightforward, one can define vectors within the framework,  
441 which includes the basis vectors in reference frames to be probabilistic. In other words, instead of using a  
442 single vector, we use a cloud of vectors given by a probability density function which we call probabilistic  
443 vectors.

444 Depending on which aspects are turned into a probabilistic vector, we can approximate the dynamics

445 with increasing fidelity: 1. A static representation as defined in the previous sections; 2. The centroid vectors  
446 are probabilistic which captures bulk location. 3. The centroid vectors are probabilistic, and the reference  
447 frames connecting modules are probabilistic which together capture bulk location and bulk orientation.  
448 4. The centroid vectors are probabilistic, and vectors to reference points are probabilistic. Captures bulk  
449 location orientation and location of reference points. For our analysis, we consider only the fourth case. It is  
450 important to note that the addition of two probabilistic vectors (p-vectors), is a convolution, while a rotation  
451 results in the rotation of the mean and covariance of the distributions.

452 To obtain these probability vectors from the conformations, we let  $M_i$  and  $\mathbf{c}_i$  be the number of points  
453 of interest and the location of the centroid for the  $i^{\text{th}}$  module. For each module, we then define a set of  
454 probability distributions given by

$$S_i = \{P(\mathbf{x}_{i,1} | \mathbf{c}_i = 0), \dots, P(\mathbf{x}_{i,M_i} | \mathbf{c}_i = 0)\} \quad (4)$$

455 where  $\mathbf{x}_{i,j}$  is the  $j^{\text{th}}$  point of interest of the  $i^{\text{th}}$  module. In addition, to  $S_i$ , we define the coupling distribution  
456  $P(\mathbf{c}_{i+1} | \mathbf{c}_i = 0)$  that describes the steady-state movement of the next module  $\mathbf{c}_{i+1}$  relative to the present one  
457  $\mathbf{c}_i$ , giving a total of  $M_i + 1$  p-vectors.

458 In order to represent these p-vectors, an appropriate probability density estimation is required that satis-  
459 fies a set of criteria: (i) it must have a parametric description for efficient storage; (ii) the density estimation  
460 must allow for rapid convolutions; (iii) the number of parameters must be “containable” so that it does not  
461 grow too large when we have many convolutions; and (iv) must be easily computable with arbitrary mo-  
462 ments. These criteria are satisfied by a Gaussian mixture model. To fit a Gaussian mixture model with the  
463 appropriate number of components, we employed a  $k$ -fold validation protocol on the conformational data.  
464 We first train a Gaussian mixture with a different number of components,  $n$ , and use the likelihood estimate  
465 of the test set from the model to score the fit. We selected the number of components that gave the highest  
466 likelihood score.

### 467 **Flexibility score to assess protein rigidity**

468 To assess the rigidity of a protein, we exploited the fact that this feature manifests itself in the model as  
469 centroids with distributions that have narrow variance. One can imagine an isosurface that expands outwards  
470 from a fixed centroid to the opposite which has the most movement. The larger the volume encapsulated  
471 by the surface the more flexible the protein and vice versa. To estimate this isosurface, for each centroid,  
472 we sampled from its Gaussian mixture density estimate giving three-dimensional points in space. We then  
473 projected these points onto a plane that is normal to the centroid backbone. Using this, we can fit a two-  
474 dimensional normal distribution to the points on the plane, from which an elliptical contour can be inferred  
475 that captures a given amount of the variance. In our case, we used the 95<sup>th</sup> percentile, which is approximately

476 two standard deviations from the mean. Connecting the ellipses with a linear interpolation gave us a pseudo-  
477 isosurface from which we then computed the volume enveloped. For convenience, we defined the cubic root  
478 of this volume as the flexibility score ( $\beta$ ) so that the score scales linearly with the number of modules.  
479 Given two constructs, the more rigid example will have a narrower envelope of movement resulting in a  
480 lower isosurface volume and  $\beta$  score. Such flexibility scores provided a convenient way to compare the  
481 rigidity of different proteins.

## 482 **Computational tools**

483 All computational simulations and analyses were run using Python version 3.11.4.

## 484 **DATA AVAILABILITY**

485 The Dynamo package used to generate all the results for this article is split into two parts. The first part,  
486 called 'dynamo', is a native Python library built in Rust for evaluating the steady-state dynamics of large  
487 bio-molecular constructs. This library is not focused on modular repeat proteins and can be used for any  
488 modular structures. It is available at: <https://github.com/seeralans/dynamo>. The second part,  
489 called 'dynamo-rp', is a Python library for coarse-grained modelling of repeat proteins and is available at:  
490 <https://github.com/seeralans/dynamo-rp>.

## 491 **ACKNOWLEDGEMENTS**

492 This work was primarily funded by UKRI grant BB/W012448/1 (to T.E.G. and F.P.) In addition, T.E.G. was  
493 supported by a Royal Society University Research Fellowship grants UF160357 and URF/R/221008, and a  
494 Turing Fellowship from The Alan Turing Institute under EPSRC grant EP/N510129/1. F.P. was supported  
495 by an EPSRC Early Career Fellowship grant EP/S017542/1. The funders had no role in study design, data  
496 collection and analysis, decision to publish or preparation of the manuscript.

## 497 **AUTHOR CONTRIBUTIONS**

498 T.E.G. and F.P. conceived the study, supervised the work, helped establish the methodology, aided with the  
499 interpretation of the results, and edited the manuscript. S.S. developed the methodology, implemented the  
500 approach, carried out all experiments and analyses, and wrote the manuscript. T.E.N. provided the initial  
501 dataset.



## 502 REFERENCES

- 503 [1] Berendsen, H. J. & Hayward, S. Collective protein dynamics in relation to function. *Current Opinion*  
504 *in Structural Biology* **10**, 165–169 (2000).
- 505 [2] Kyte, J. Molecular considerations relevant to the mechanism of active transport. *Nature* **292**, 201–204  
506 (1981).
- 507 [3] Kolasangiani, R., Bidone, T. C. & Schwartz, M. A. Integrin conformational dynamics and mechan-  
508 otransduction. *Cells* **11**, 3584 (2022).
- 509 [4] Dou, J. *et al.* De novo design of a fluorescence-activating  $\beta$ -barrel. *Nature* **561**, 485–491 (2018).
- 510 [5] Huang, P.-S. *et al.* De novo design of a four-fold symmetric TIM-barrel protein with atomic-level  
511 accuracy. *Nat Chem Biol* **12**, 29–34 (2016).
- 512 [6] Silva, D.-A. *et al.* De novo design of potent and selective mimics of IL-2 and IL-15. *Nature* **565**,  
513 186–191 (2019).
- 514 [7] Röthlisberger, D. *et al.* Kemp elimination catalysts by computational enzyme design. *Nature* **453**,  
515 190–195 (2008).
- 516 [8] Langan, R. A. *et al.* De novo design of bioactive protein switches. *Nature* **572**, 205–210 (2019).
- 517 [9] Tokuriki, N. & Tawfik, D. S. Protein dynamism and evolvability. *Science* **324**, 203–207 (2009).
- 518 [10] Mori, T., Miyashita, N., Im, W., Feig, M. & Sugita, Y. Molecular dynamics simulations of biological  
519 membranes and membrane proteins using enhanced conformational sampling algorithms. *Biochimica*  
520 *et Biophysica Acta (BBA) - Biomembranes* **1858**, 1635–1651 (2016).
- 521 [11] Harpole, T. J. & Delemotte, L. Conformational landscapes of membrane proteins delineated by en-  
522 hanced sampling molecular dynamics simulations. *Biochimica et Biophysica Acta (BBA) - Biomem-*  
523 *branes* **1860**, 909–926 (2018).
- 524 [12] Loschwitz, J., Olubiyi, O. O., Hub, J. S., Strodel, B. & Poojari, C. S. Computer simulations of  
525 protein–membrane systems. In Strodel, B. & Barz, B. (eds.) *Progress in Molecular Biology and*  
526 *Translational Science*, vol. 170 of *Computational Approaches for Understanding Dynamical Systems:*  
527 *Protein Folding and Assembly*, 273–403 (Academic Press, 2020).
- 528 [13] Orellana, L. Large-scale conformational changes and protein function: Breaking the in silico barrier.  
529 *Frontiers in Molecular Biosciences* **6** (2019).
- 530 [14] Kaynak, B. T. *et al.* Sampling of protein conformational space using hybrid simulations: A critical  
531 assessment of recent methods. *Frontiers in Molecular Biosciences* **9** (2022).
- 532 [15] Park, J.-K., Jernigan, R. & Wu, Z. Coarse Grained Normal Mode Analysis vs. Refined Gaussian  
533 Network Model for protein residue-level structural fluctuations. *Bull Math Biol* **75**, 124–160 (2013).

- 534 [16] Kmiecik, S. *et al.* Coarse-grained protein models and their applications. *Chem. Rev.* **116**, 7898–7936  
535 (2016).
- 536 [17] Joshi, S. Y. & Deshmukh, S. A. A review of advancements in coarse-grained molecular dynamics  
537 simulations. *Molecular Simulation* **47**, 786–803 (2021).
- 538 [18] Clementi, C. Coarse-grained models of protein folding: Toy models or predictive tools? *Current*  
539 *Opinion in Structural Biology* **18**, 10–15 (2008).
- 540 [19] Tozzini, V. Coarse-grained models for proteins. *Current Opinion in Structural Biology* **15**, 144–150  
541 (2005).
- 542 [20] Atilgan, A. R. *et al.* Anisotropy of fluctuation dynamics of proteins with an elastic network model.  
543 *Biophysical Journal* **80**, 505–515 (2001).
- 544 [21] Fuglebakk, E., Reuter, N. & Hinsen, K. Evaluation of protein elastic network models based on an  
545 analysis of collective motions. *J. Chem. Theory Comput.* **9**, 5618–5628 (2013).
- 546 [22] López-Blanco, J. R. & Chacón, P. New generation of elastic network models. *Current Opinion in*  
547 *Structural Biology* **37**, 46–53 (2016).
- 548 [23] Dehouck, Y. & Bastolla, U. Why are large conformational changes well described by harmonic normal  
549 modes? *Biophysical Journal* **120**, 5343–5354 (2021).
- 550 [24] Yang, L., Song, G. & Jernigan, R. L. How well can we understand large-scale protein motions using  
551 normal modes of elastic network models? *Biophysical Journal* **93**, 920–929 (2007).
- 552 [25] Ventura, C., Banerjee, A., Zacharopoulou, M., Itzhaki, L. S. & Bahar, I. Tandem-repeat proteins con-  
553 formational mechanics are optimized to facilitate functional interactions and complexations. *Current*  
554 *Opinion in Structural Biology* **84**, 102744 (2024).
- 555 [26] M. Delucchi, M., Schaper, E., Sachenkova, O., Elofsson, A. & Anisimova, M. A new census of protein  
556 tandem repeats and their relationship with intrinsic disorder. *Genes* **11**, 407 (2020).
- 557 [27] Paladin, L. *et al.* Repeatsdb in 2021: improved data and extended classification for protein tandem  
558 repeat structures. *Nucleic Acids Research* **49**, D452–D457 (2021).
- 559 [28] Kobe, B. & Kajava, A. V. The leucine-rich repeat as a protein recognition motif. *Current Opinion in*  
560 *Structural Biology* **11**, 725–732 (2001).
- 561 [29] Li, J., Mahajan, A. & Tsai, M.-D. Ankyrin repeat: A unique motif mediating protein-protein interac-  
562 tions. *Biochemistry* **45**, 15168–15178 (2006).
- 563 [30] Grove, T. Z., Cortajarena, A. L. & Regan, L. Ligand binding by repeat proteins: Natural and designed.  
564 *Current Opinion in Structural Biology* **18**, 507–515 (2008).
- 565 [31] Parmeggiani, F. & Huang, P.-S. Designing repeat proteins: A modular approach to protein design.  
566 *Current Opinion in Structural Biology* **45**, 116–123 (2017).

- 567 [32] Park, K. *et al.* Control of repeat-protein curvature by computational protein design. *Nat Struct Mol*  
568 *Biol* **22**, 167–174 (2015).
- 569 [33] Ferreiro, D. U., Walczak, A. M., Komives, E. A. & Wolynes, P. G. The energy landscapes of repeat-  
570 containing proteins: Topology, cooperativity, and the folding funnels of one-dimensional architectures.  
571 *PLoS Comput Biol* **4**, e1000070 (2008).
- 572 [34] Espada, R. *et al.* Repeat proteins challenge the concept of structural domains. *Biochemical Society*  
573 *Transactions* **43**, 844–849 (2015).
- 574 [35] Kaynak, B. T. *et al.* Cooperative mechanics of pr65 scaffold underlies the allosteric regulation of the  
575 phosphatase pp2a. *Structure* **31**, 607–618.e3 (2023).
- 576 [36] Synakewicz, M. *et al.* Unraveling the mechanics of a repeat-protein nanospring: From folding of  
577 individual repeats to fluctuations of the superhelix. *ACS Nano* **16**, 3895–3905 (2022).
- 578 [37] Yeh, C.-T., Brunette, T., Baker, D., McIntosh-Smith, S. & Parmeggiani, F. Elfin: An algorithm for  
579 the computational design of custom three-dimensional structures from modular repeat protein building  
580 blocks. *Journal of Structural Biology* **201**, 100–107 (2018).
- 581 [38] Brunette, T. J. *et al.* Exploring the repeat protein universe through computational protein design.  
582 *Nature* **528**, 580–584 (2015).
- 583 [39] Brunette, T. J. *et al.* Modular repeat protein sculpting using rigid helical junctions. *Proceedings of the*  
584 *National Academy of Sciences* **117**, 8870–8875 (2020).
- 585 [40] Tyka, M. D. *et al.* Alternate states of proteins revealed by detailed energy landscape mapping. *Journal*  
586 *of Molecular Biology* **405**, 607–618 (2011).
- 587 [41] Shang, Z. *et al.* High-resolution structures of kinesin on microtubules provide a basis for nucleotide-  
588 gated force-generation. *eLife* **3**, e04686 (2014).
- 589 [42] Kerns, S. J. *et al.* The energy landscape of adenylate kinase during catalysis. *Nat Struct Mol Biol* **22**,  
590 124–131 (2015).
- 591 [43] Praetorius, F. *et al.* Design of stimulus-responsive two-state hinge proteins. *Science* **381**, 754–760  
592 (2023).
- 593 [44] Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589  
594 (2021).
- 595 [45] Zheng, L.-E., Barethiya, S., Nordquist, E. & Chen, J. Machine learning generation of dynamic protein  
596 conformational ensembles. *Molecules* **28** (2023).
- 597 [46] Strokach, A. & Kim, P. M. Deep generative modeling for protein design. *Current Opinion in Structural*  
598 *Biology* **72**, 226–236 (2022).
- 599 [47] Noé, F., De Fabritiis, G. & Clementi, C. Machine learning for protein folding and dynamics. *Current*

600 *Opinion in Structural Biology* **60**, 77–84 (2020).

601 [48] Audagnotto, M. *et al.* Machine learning/molecular dynamic protein structure prediction approach to  
602 investigate the protein conformational ensemble. *Scientific Reports* **12**, 10018 (2022).

603 [49] Janson, G., Valdes-Garcia, G., Heo, L. & Feig, M. Direct generation of protein conformational ensem-  
604 bles via machine learning. *Nature Communications* **14**, 774 (2023).

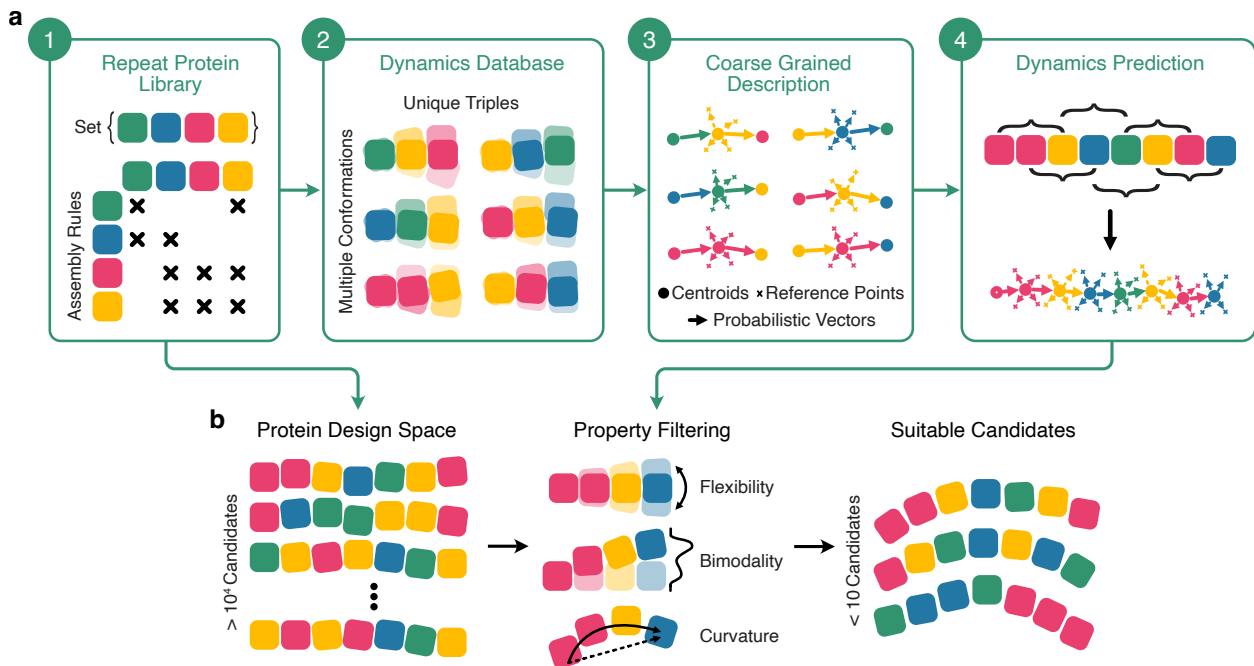
605 [50] Heinig, M. & Frishman, D. STRIDE: a web server for secondary structure assignment from known  
606 atomic coordinates of proteins. *Nucleic Acids Research* **32**, W500–W502 (2004).

607 [51] Kabsch, W. & Sander, C. Dictionary of protein secondary structure: Pattern recognition of hydrogen-  
608 bonded and geometrical features. *Biopolymers* **22**, 2577–2637 (1983).

609 [52] Nivón, L. G., Moretti, R. & Baker, D. A pareto-optimal refinement method for protein design scaffolds.  
610 *PLoS ONE* **8**, e59004 (2013).

611 [53] Conway, P., Tyka, M. D., DiMaio, F., Konerding, D. E. & Baker, D. Relaxation of backbone bond  
612 geometry improves protein energy landscape modeling. *Protein Science* **23**, 47–55 (2014).

613 **FIGURES AND CAPTIONS**



614

615 **Figure 1: Overview of the coarse-grained model for the structural and dynamical design of**

616 **modular proteins.** (a) Workflow for generating a coarse-grained model of protein structural dynamics.

617 The model relies on an underlying repeat protein library that contains well defined protein modules and

618 rules for how these can be assembled (Step 1). Using this information, a dynamics database of all three

619 module proteins capturing the relative movements between all atoms in the protein is generated (Step 2).

620 This can be via detailed simulations (e.g., using Rosetta) or derived from experimental data. The resultant

621 dynamics database is then used to build coarse-grained descriptions of each triplet of modules, defining

622 centroids for each module and an arbitrary number of reference points, e.g., the end points of alpha helices

623 (Step 3). Finally, these coarse-grained descriptions are stitched together to enable the efficient propagation

624 of movements of larger proteins built from the repeat protein library (Step 4). (b) The coarse-grained

625 model can be used to accelerate the design of proteins built using the repeat protein library. A typical library

626 defines a vast potential design space that would be impossible to exhaustively search. The speed of the

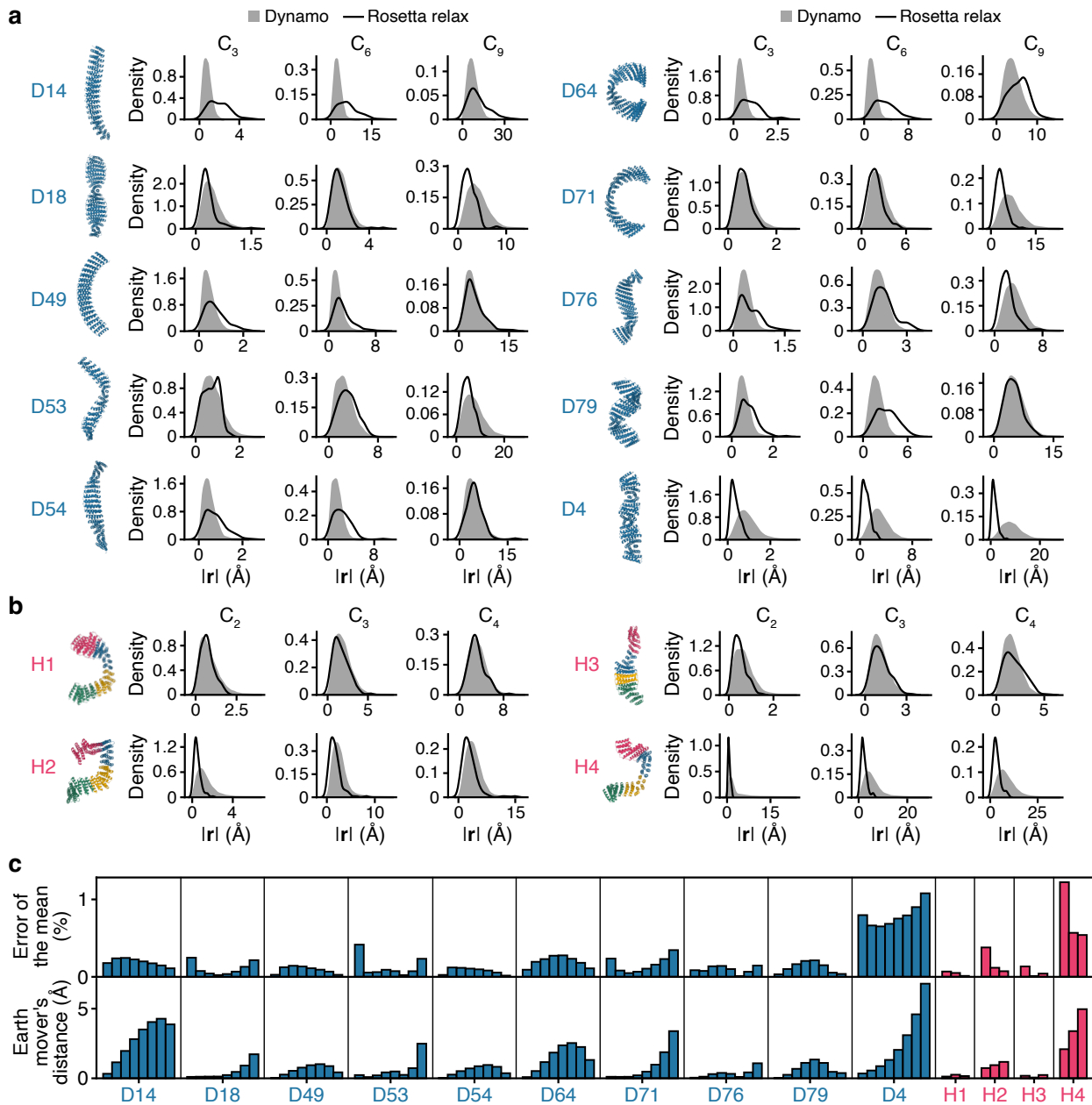
627 coarse-grained model allows for large regions of this space to be probed (e.g., millions of designs) and key

628 structural and dynamical properties measured (e.g., flexibility or potential for multi-state dynamics). This

629 information can be used to guide future areas to explore and the properties of the simulated designs can be

630 filtered on properties that are essential to the desired function of the protein. Using this approach a targeted

631 set of designs that can be feasibly built is output for detailed experimental testing.



632

633

634

635

636

637

638

639

640

641

642

643

**Figure 2: Comparison of conformational data obtained using the Rosetta relax protocol**

**and our model (Dynamo).** (a) Distributions capturing the fluctuations in the position of centroids ( $C_x$ )

relative to their respective mean. Each centroid has a mean position, and the displacement from this mean

and samples of centroid positions obtained from our model (grey filled distribution) or through Rosetta relax

(solid black line) is given by  $\mathbf{r}$ . Distributions shown for ten homogeneous 9-module constructs. A molecular

visualisation of the construct is shown on the right of each plot. (b) Similar distributions as described in

panel (a) for four heterogeneous 4-module constructs. The modules in the heterogeneous constructs are:

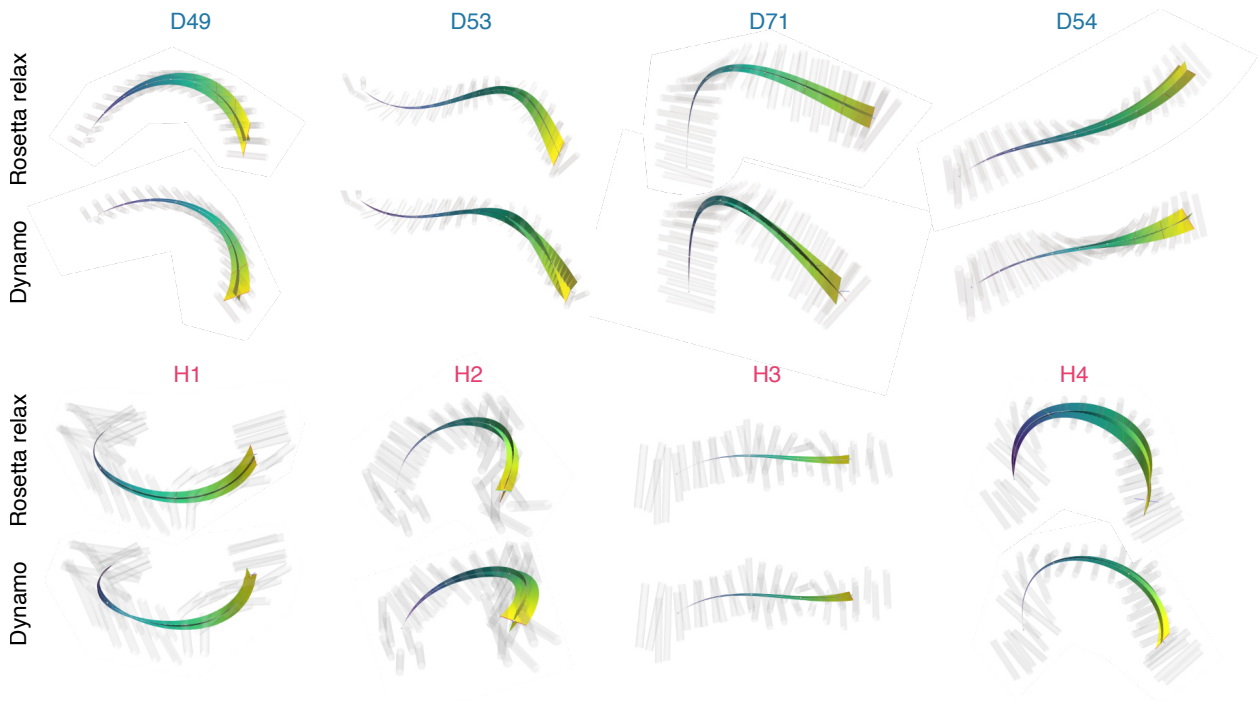
H1 = D14\_j1\_D14x4\_188; H2 = D14\_j1\_D18x4\_31; H3 = D49\_j1\_D49x4\_49; H4 = D79\_j1\_D14x4\_131. (c)

Comparison between the probability densities in panels (a) and (b) between the model and the Rosetta relax

data for each centroid: (top) percentage error of the mean of the distributions, (bottom) a full distribution

comparison using the earth mover's distance. Bars correspond to centroids  $C_2$  to  $C_9$  (left to right) for the

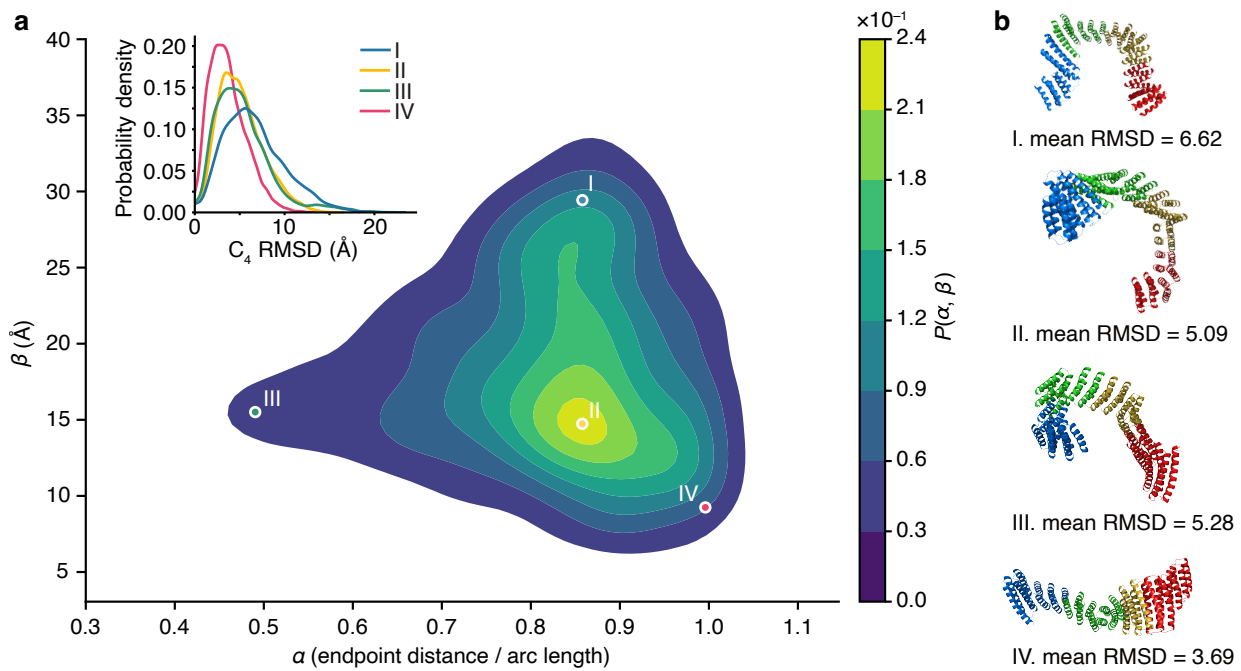
644 9-module homogeneous constructs and  $C_2$  to  $C_4$  (left to right) for the 4-module heterogeneous constructs.



645

646 **Figure 3: Visualisation of the structural dynamics of several modular protein designs.** Data  
647 shown for the Rosetta relax protocol (top) compared to our model (bottom). The two fins are parallel to the  
648 two directions with the most movement and envelope the 95<sup>th</sup> percentile of the centroid density distributions,  
649 while the colour and width of the fins correspond to the largest movement (dark blue to yellow denoting  
650 small to large movements, respectively). The alpha helices are represented by grey semi-transparent cylin-  
651 ders. The modules in the heterogeneous constructs are: H1 = D14\_j1\_D14x4\_188; H2 = D14\_j1\_D18x4\_31;  
652 H3 = D49\_j1\_D49x4\_49; H4 = D79\_j1\_D14x4\_131.





653

654 **Figure 4: Exploring the design landscape of all 4-module protein chains containing 2,978**

655 **unique designs.** (a) Smoothed density plot of all 4-module protein chains from our model. The position

656 of four selected designs is indicated. Note that  $0 \leq \alpha \leq 1$ , but kernel density estimates can give non zero

657 probability outside of this range, for clarity of visualisation we have chosen not to remove the small region of

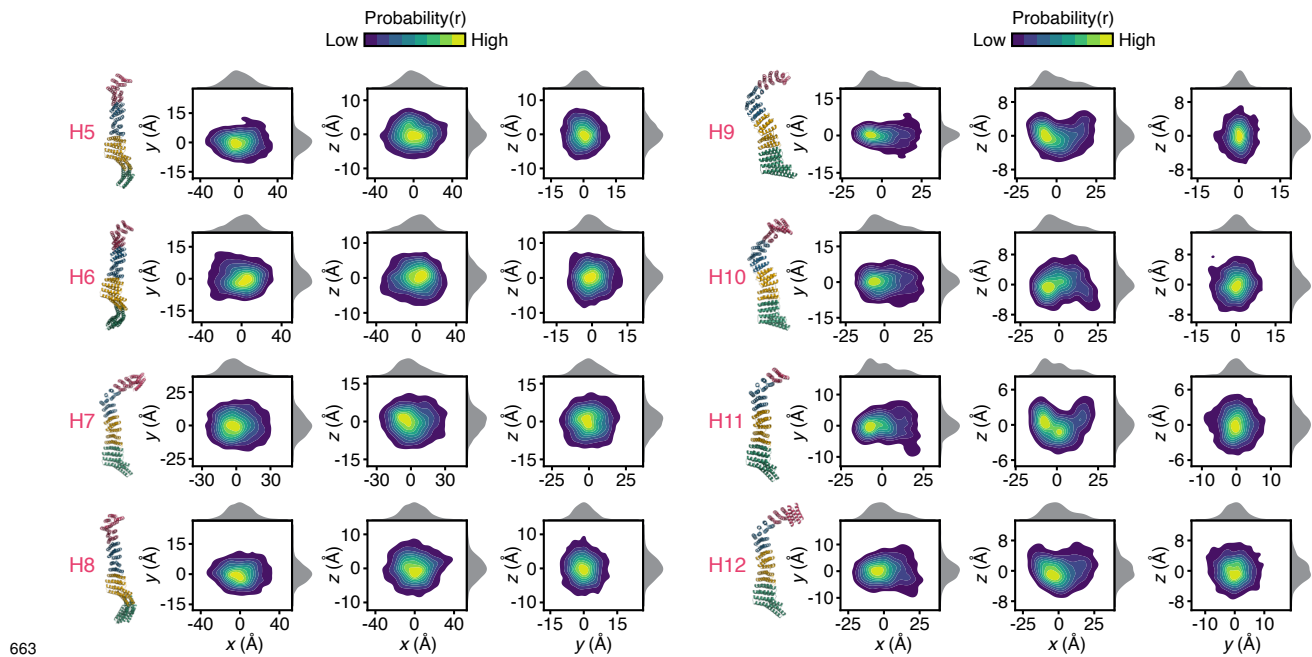
658 probability estimate for  $\alpha > 1$ . Insert shows the probability distribution of the root mean squared deviation

659 (RMSD) of the final module in the chain ( $C_4$ ) for the four highlighted designs (I–IV), as calculated from full

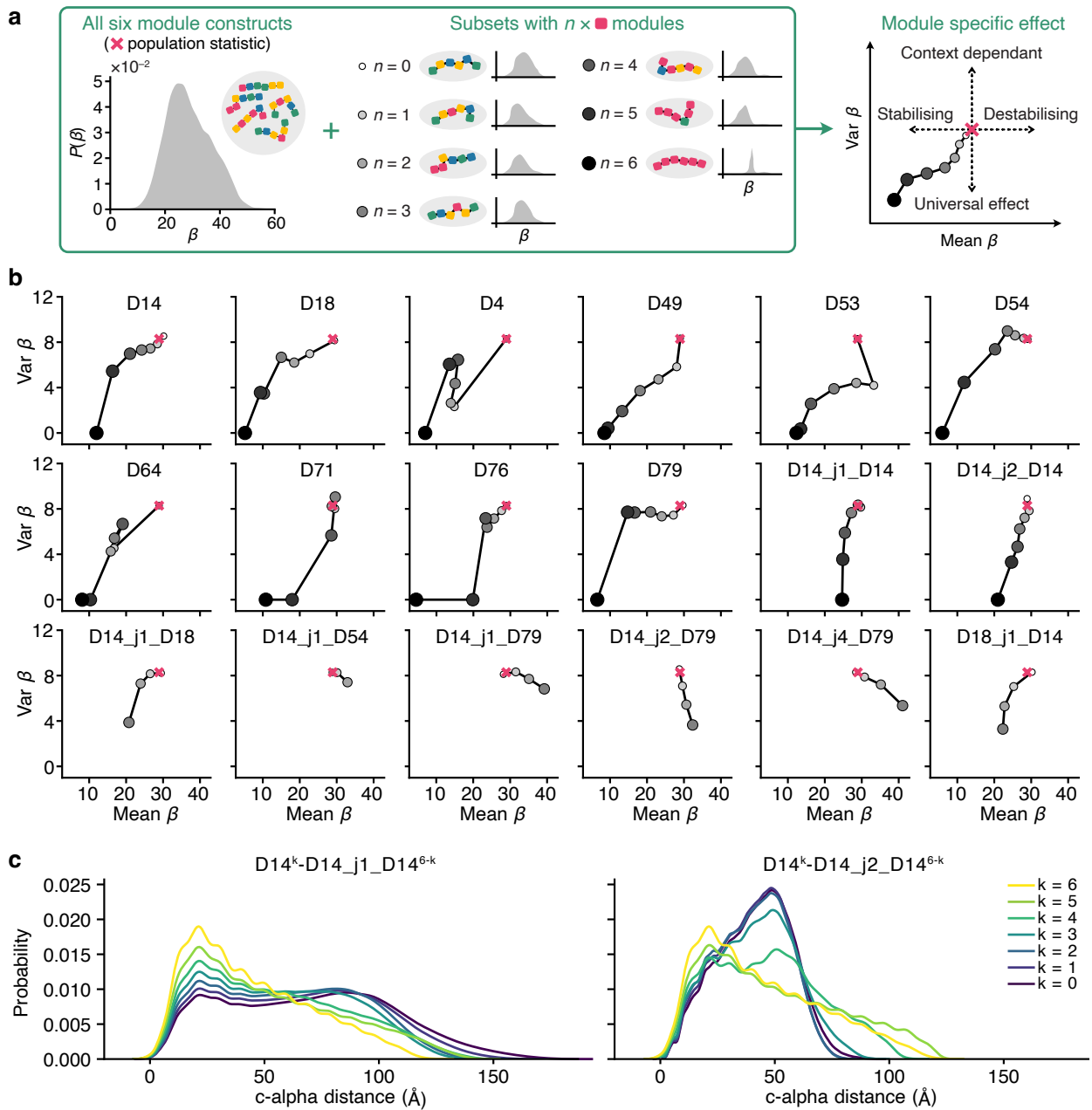
660 length Rosetta relax runs. This distribution captures the general range of dynamic movements the module

661 experiences and correlates with the flexibility ( $\beta$ ). (b) Structural visualisations of specific protein designs

662 highlighted in panel a. Individual modules denoted in different colours from N- (blue) to C-terminus (red).

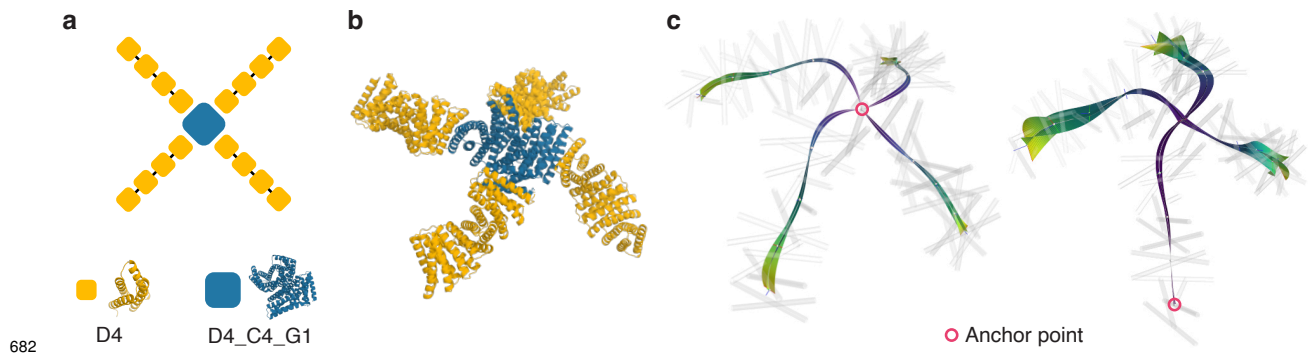


664 **Figure 5: Exploring potential multi-stable behaviour of 4-module constructs.** The eight most  
665 highly-ranked constructs after scoring each on their potential for multi-stable behaviour. The position of the  
666 centroid of the last module,  $r$ , is visualised as a probability distribution projected onto orthogonal planes.  
667 Each distribution is built by sampling  $10^6$  centroid positions from the model. The modules in the heteroge-  
668 nous constructs are: H5 = D14\_j1\_D14x4\_239; H6 = D14\_j1\_D14x4\_241; H7 = D79\_j2\_D14x4\_117; H8  
669 = D14\_j1\_D14x4\_240; H9 = D18\_j1\_D14x4\_117; H10 = D14\_j1\_D14x4\_117; H11 = D14x4\_117; H12 =  
670 D49\_j1\_D14x4\_117.



671

672 **Figure 6: Analysing the role of each module on the flexibility of 6-module constructs.** (a) The  
 673 module specific effect is summarised by a plot showing the mean and variance of  $\beta$  probability distributions,  
 674 with points for constructs containing 0 to 6 instances of the modules of interest (small white filled to large  
 675 black filled circle). A summary statistic is plotted (red cross) of the entire set of all 6-module designs.  
 676 Changes in the mean and variance of  $\beta$  as the number of instances of a module increase relate to stabilising or  
 677 destabilising effects that are either context dependant or universal for all other modules. (b) Module specific  
 678 effect plots for each module in every 6-module construct. (c) Probability density of all-to-all pairwise  
 679 distances of alpha carbons in a construct containing a combination of D14 and D14\_j1\_D14 or D14\_j2\_D14,  
 680 normalised to the maximum distance of a pure D14 construct (177 Å).  $k$  indicates the number of D14  
 681 modules in the protein.



682  
683 **Figure 7: Estimating dynamics of a multi-chain tetramer** (a) A schematic representation of a  
684 multi-chain construct and the two distinct modules used in the construct. (b) Molecular visualization of the  
685 multi-chain construct. We use D4\_C4\_G1 hub with four D4 modules connected to each of the four arms of  
686 the hub. (c) A visualization of the structural dynamics predicted using our model with two different anchor  
687 points (red circle).