
Normal Mode Diffusion: Towards Dynamics-Informed Protein Design

Anonymous Authors¹

Abstract

We introduce a new method that uses normal mode analysis (NMA) to condition diffusion models for protein design to create proteins with specific dynamical properties – that is, their lowest non-trivial normal mode moves a selected set of residues in a targeted way. We demonstrate that our approach is feasible by conditioning an equivariant graph-diffusion model for protein backbone generation to design molecules with a pre-defined lowest normal mode. Our work represents a first step towards incorporating dynamical behaviour in protein design, and may open the door to designing more flexible and functional proteins in the future.

1. Introduction

Generative artificial intelligence (AI) has demonstrated significant advances in molecular design. In protein design, denoising diffusion models allow the controllable generation of a range of proteins (Watson et al., 2022; Ingraham et al., 2022). Current models allow researchers to condition protein designs on desirable characteristics such as shape, scaffolding functional motifs, and other sequence and structural properties.

However, despite the extensive range of properties that can be conditioned on, no methods to condition generative models on protein dynamical properties have been introduced. Here, we present a novel method to condition diffusion models on protein dynamics data derived from normal mode analysis (NMA) (Bahar et al., 2010). Our approach expands the potential for generating proteins with specific dynamical properties, crucial for functional applications.

NMA offers unique advantages for generative AI-driven protein design. First, it presents a fast route to a rough

hypothesis of *dynamics* for any protein purely based on structure, without the need for expensive molecular dynamics information (Bahar et al., 2010). Second, despite the fact that NMA is less sophisticated than molecular dynamics simulations, the lowest (5-15) non-trivial normal modes have been shown to capture a substantial portion of functional motions in proteins (Bahar et al., 2010). Prominent examples include dihydrofolate reductase (DHFR) (Bahar et al., 1997), lysozyme (Gibrat & Gō, 1990), and adenylate kinase (AdK) (Tama & Sanejouand, 2001), among others.

Our results show how protein backbone generation can be conditioned on such that their lowest normal mode moves a select set of residues move in a targeted way. This provides a route to designing proteins with tailored dynamical properties which may open new avenues for protein engineering and drug discovery. Our contributions are:

1. To our knowledge, this is the first proposal to introduce explicit conditioning on dynamical properties into diffusion models for protein design.
2. We evaluate our method in a proof-of-concept setting on its ability to generate protein structures that are (1) realistic, (2) novel and (3) follow the specified dynamics condition. Our experiments demonstrate that conditioning on NMA-dynamics is viable and can be incorporated in current diffusion models.

2. Background and related work

Diffusion models Denoising diffusion probabilistic models (DDPM) (Sohl-Dickstein et al., 2015; Ho et al., 2020) have been applied in a variety of tasks within and outside biology, such as image synthesis (Dhariwal & Nichol, 2021; Kong et al., 2021), drug design (Schneuring et al., 2023) and protein design (Watson et al., 2022; Ingraham et al., 2022). Given a data sample x_0 , DDPMs learn are trained to remove random noise that is added to the sample. By doing so, they learn to reconstruct data samples from noise. We refer the reader to Ho et al. (2020) for a good introduction.

Normal mode analysis (NMA) NMA is a technique to investigate the vibrational dynamics of a system of molecules (Bahar et al., 2010). It is based on a harmonic approximation of the potential energy around a given protein structure.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

From this approximation, the vibrational frequencies and corresponding directions of oscillation (normal modes) arise as eigenvectors and eigenvalues of the potential energy Hessian matrix. Compared to molecular dynamics simulations, NMA is computationally efficient and can often distill complex dynamics into a few dominant modes (more in App. B).

3. Methods

Setup We consider the simplified setting of monomeric protein chains which we represent with only the most relevant features as being an *ordered* point cloud $x_0 \in \mathbb{R}^{n \times 3}$ constituted only of the N-to-C ordered list of alpha-carbon coordinates of its residues, without any side-chain information. As training and validation data we extract 10'000 high-resolution ($< 3 \text{ \AA}$) samples from CATHv4.3 (Orengo et al., 1997) of lengths between 25-100 amino acids (Details in App. C).

Diffusion loop, denoiser and training We follow Hoogeboom et al. (2022) and use the Markovian DDPM formulation of Ho et al. (2020) for our forward process. We use the noise schedule from Hoogeboom et al. (2022) with 500 steps and subtract the center of mass from the noise for an equivariant diffusion process (Hoogeboom et al., 2022).

As denoising model $\epsilon_\theta(x_t, t)$, we use the Geometric Vector Perceptron (GVP-GNN) (Jing et al., 2021), where the normalised time-step t is added as node-feature, and the initial layer norms in the embeddings are dropped. We perform message passing on a fully connected C_α graph with positional encoding to capture the chain structure. The denoiser is trained using the denoising objective from Ho et al. (2020) for 500 epochs with learning rate 1e-4.

Unconditional sampling We follow (Hoogeboom et al., 2022), and start by sampling $x_T \sim \mathcal{N}(0, 1)$, subtracting the center of mass, and applying the denoising steps using

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\sqrt{1-\alpha_t}}{1-\bar{\alpha}_t} \epsilon_\theta(x_t, t) \right) + \eta(1-\alpha_t)z, \quad (1)$$

with $z \sim \mathcal{N}(0, 1)$ with subtracted center of mass. Motivated by Song et al. (2020), we introduce an empirical noise scale η , with $\eta = 1$ being the DDPM formulation in Ho et al. (2020). We find that the quality of the unconditioned samples is generally improved for small η . We therefore use a deterministic reverse process ($\eta = 0$) during all experiments below.

NMA-conditioned sampling We consider the following situation: Given a set \mathcal{C} of residues of interest – for example a functional motif – we would like to specify their dominant relative motion.

To specify the relative motion, we define a *target matrix*

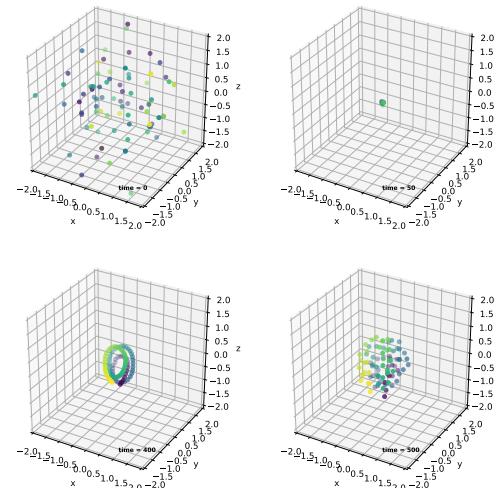


Figure 1: Snapshots of the reverse process for unconditional backbone generation. Point-cloud samples from a Gaussian prior with center-of-mass zero are progressively denoised. We observe that the reverse process firstly aligns the backbone residues in a chain, before expanding the structure to biological sizes. Structural details such as alpha helices or beta strands are formed towards the end of the process.

$y \in \mathbb{R}^{|\mathcal{C}| \times 3}$ of normal mode component vectors $y_i \in \mathbb{R}^3$ for each residue i in \mathcal{C} . The target matrix may be manually specified, or extracted from the normal modes of a functional motif of interest in a target protein. In this work, we wanted to evaluate against diverse motions and therefore choose samples according to a *strain-energy* calculation: We randomly sample a target protein in the holdout data, perform NMA and identify the most flexible part of the protein by calculating the strain energy of each node (Hinsen & Kneller, 1999). We then choose the lowest normal mode component of the 10 consecutive nodes for which the summed energy is largest as our target matrix y .

Conditioning with gradient descent. We take inspiration from classifier-based guidance (Dhariwal & Nichol, 2021), where the score of the data is conditioned on the target y by Bayes rule:

$$\nabla_{x_t} \ln p(x_t|y) = \nabla_{x_t} \ln p(y|x_t) + \nabla_{x_t} \ln p(x_t) \quad (2)$$

The key idea is that instead of training a model to predict $p(y|x_t)$, we exploit that coarse-grained NMA is robust to slight variations in structure (Bahar et al., 2010) and can be computed differentiably. Therefore, once a rough structural hypothesis is formed we can use gradient descent on an NMA based loss to induce a probability flow towards samples that satisfy the condition y :

$$\Delta x = -\gamma(x_t, t) \nabla_{x_t} l(y, v(x_t)) \quad (3)$$

110 Here $v(x_t)$ is the lowest non-trivial normal mode of x_t ,
 111 $\gamma(x_t, t)$ is a guidance scale and $l(y, v(x_t))$ is the loss. To
 112 ensure meaningful NMA results, we require a somewhat
 113 protein-like structure x_t . We therefore start conditioning
 114 only from $t < t_{\text{start}}$ onward. As seen in Figure 1, a mean-
 115 ingful chain structure emerges in the middle of our reverse
 116 process and the physical distances between residues are only
 117 recovered towards the last quarter of the diffusion process.
 118 Based on these observations, we found a $t_{\text{start}} = T/5$ to
 119 work well for our model.

120 Upon sampling from the reverse process and starting from
 121 $t < t_{\text{start}}$, we condition in the following way. We (1) inflate
 122 the current structure to physical size such that mean chain
 123 distance is 3.8 Å, (2) extract the lowest non-trivial mode of
 124 the structure using Hinsen force-field (Hinsen & Kneller,
 125 1999) parametrisation with a 16 Å cut-off, (3) subset the *current mode*
 126 to the mode components for conditioned residues
 127 $v(x_t) \in \mathbb{R}^{|\mathcal{C}| \times 3}$.

128 The structure x_t is then optimised towards the target y by
 129 computing the loss (Eq. 4) and performing gradient descent.
 130 In practice we found it beneficial to split each conditioning
 131 time step into $r = 5$ smaller steps, interleaving a partial de-
 132 noising (with rescaled amplitudes) with a gradient descent
 133 step with recalculated loss. To balance the magnitudes of
 134 the denoising and the conditioning update, we set $\gamma(x_t, t)$
 135 to the maximal magnitude of the denoising update.

136 The conditioning loss (Equation 4) is chosen as a simple
 137 combination of amplitude and angle terms between all pair-
 138 wise residues.

$$139 l_{\text{NMA}}(y, v(x_t)) = l_{\text{angle}}(y, v(x_t)) + 2l_{\text{ampl}}(y, v(x_t)) \quad (4)$$

$$140 l_{\text{angle}} = \sum_{i,j \in \mathcal{C}} |\cos(y_i, y_j) - \cos(v(x_t)_i, v(x_t)_j)| \quad (5)$$

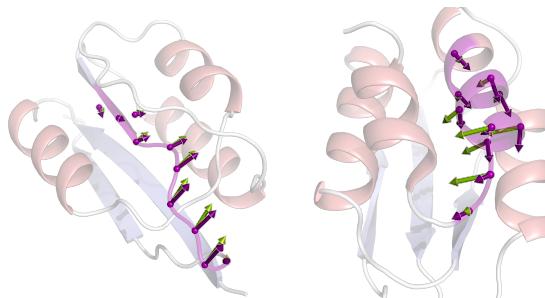
$$141 l_{\text{ampl}} = \sum_{i \in \mathcal{C}} \left| \frac{\|y_i\|}{\|y\|} - \frac{\|v(x_t)_i\|}{\|v(x_t)\|} \right| \quad (6)$$

142 This choice extracts invariants from the target matrix y
 143 which are independent of the reference frame. Further, it
 144 ensures rotational equivariance of the gradient with respect
 145 to rotations of x_t . The amplitude terms are normalised such
 146 that only their relative sizes matter, consistent with the fact
 147 that amplitude information from NMA can only make rela-
 148 tive statements about the participation of a given residue in a
 149 mode (Bahar et al., 2010). For the combined loss, the l_{ampl} is
 150 scaled by 2, such that its contribution is similar in magnitude
 151 to l_{angle} . This simple loss does not consider higher order
 152 correlations (e.g. among the motion of triplets of residues),
 153 but could readily be extended to do so.

154 4. Results

155 Our goal is to evaluate the efficacy of NMA Diffusion in
 156 generating realistic backbones with specific dynamic prop-
 157 erties. This prompts us to address three central questions:
 158 (1) Do the generated samples represent realistic proteins?
 159 (2) If they represent realistic proteins, do they display the
 160 desired normal mode dynamics? (3) Are the samples novel
 161 compared to the training set?

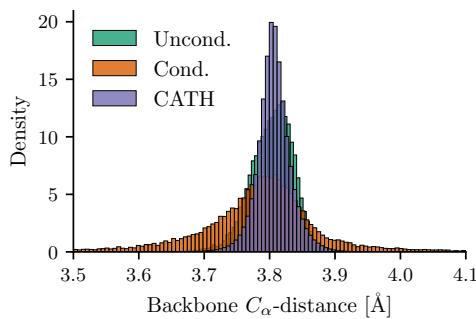
162 To assess these questions, we use 300 hold-out target struc-
 163 tures of various lengths from our CATH dataset and extract
 164 an NMA target condition as explained above. We then gen-
 165 erate 3 unconditioned and 4 NMA-conditioned samples for
 166 each target condition. We then filter the samples for realis-
 167 tic backbones, requiring that the mean C_α distance along
 168 the backbone is within 0.05 Å of the average C_α distance
 169 of 3.8 Å (Voet & Voet, 2010). This reduces the number of
 170 samples by approximately 25% and 50% for unconditional
 171 and conditional samples respectively. Of these, we select
 172 the sample with the lowest NMA loss (Eq. 4), leaving with
 173 approximately 300 samples per sampling procedure. Two
 174 samples are shown in Figure 2 (more in App. F).



175 **Figure 2: Left:** conditioned sample with NMA-loss $l_{\text{NMA}} =$
 176 0.11. **Right:** not conditioned sample with $l_{\text{NMA}} = 0.51$.
 177 Purple arrows represent the target, green displacements in
 178 the novel protein.

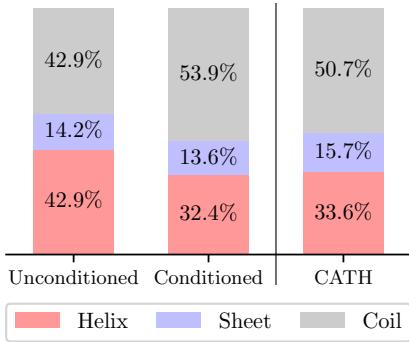
179 **Filtered samples represent realistic proteins** For the
 180 filtered samples, we compute the C_α backbone distances
 181 (Fig. 3). The unconditioned samples follow the C_α
 182 distance distribution of the training set well, indicating that
 183 the denoiser can generate realistic backbones. The NMA-
 184 conditioned samples show a similar distribution, but have
 185 heavy tails extending to occasionally unrealistic distances.
 186 When filtering samples with unrealistic distances, we lose
 187 another 3x of samples compared to the unconditioned case.
 188 We attribute the distortion of backbone distances to the
 189 NMA-loss, which does not actively promote realistic back-
 190 bone distances. Currently, the backbone distances are only
 191 corrected by the denoiser, which for our simple denoiser ap-
 192 pears to be insufficient. Since we did not spend much time
 193 in tuning and training our denoiser for this work, we believe
 194 that a more sophisticated denoiser (Ingraham et al., 2022)
 195 or an extra term in the conditioning to encourage maintain-
 196 ing realistic backbone distances can fix the drop in sample
 197 quality. Overall, we conclude that the NMA-conditioned

165 sampling can generate realistic backbones, but at the cost of
 166 requiring more samples to find a realistic backbone.
 167



179 Figure 3: C_α backbone distances of the generated samples
 180 compared to the training set.

182 We further compute secondary structure features from the
 183 C_α positions with the P-SEA algorithm (Labesse et al.,
 184 1997) and compare it to the distribution of true proteins
 185 in our dataset (Fig. 4). The unconditioned samples have a
 186 similar secondary structure distribution as the training set,
 187 with a slight overrepresentation of alpha helices. This is
 188 in line with previous work (Watson et al., 2022) Interest-
 189 ingly, the NMA-conditioned sampling remedies this bias
 190 and reproduces the secondary structure distribution of the
 191 training distribution better. This shows that the NMA con-
 192 ditioning does destroy the secondary structure and indeed
 193 may allow correcting for the bias of unconditioned sampling.
 194 When looking at generated samples (App. F) we find that
 195 the generated samples still often show slightly unrealistic
 196 packing, but that there are also samples which are plausible
 197 combination of structural motifs (e.g. Fig. 2).



211 Figure 4: Secondary structure distribution of the generated
 212 samples compared to the training set.

214 **Conditioned samples exhibit the targeted normal mode**
 215 To examine whether conditioning the selected residues to
 216 have the targeted lowest normal mode component, we com-
 217 pute the NMA loss (Eq. 4) for the realistic, generated sam-
 218 ples and compare it to the unconditioned samples (Fig. 5).
 219

We see that the NMA-conditioned samples are significantly enriched towards low loss values compared to the unconditional samples. In the examples in Fig. 2 and App. F we also see that lowest normal mode of the generated samples (green) is in good agreement with the target (purple). This indicates that NMA-conditioning is effective in generating samples with the targeted lowest normal mode component.

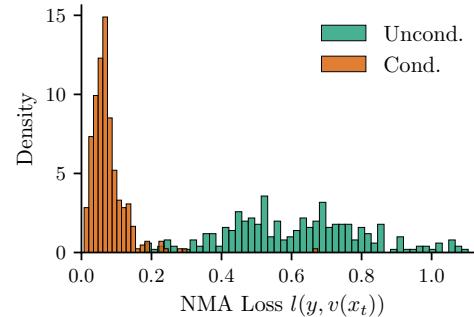


Figure 5: Histograms of the NMA-loss (Eq. 4) for conditioned and not conditioned samples.

Samples are novel. Finally, to ensure that the model did not overfit we compute the TM-score (Zhang & Skolnick, 2005) between the generated samples and the best matching target fragment of the same size in the training set (App. E). Both, conditioned and unconditioned samples, have TM-scores of about 0.3-0.4, indicating that they are novel compared to the training set and do not simply memorise the training data. Importantly, NMA-conditioning does not have a negative effect on the novelty of the samples.

5. Conclusion and further work

We introduced a novel framework to condition diffusion models for protein design on protein dynamics information with normal mode analysis. Our analysis demonstrates that with our *NMA-Diffusion* sampling procedure, it is possible to condition a generative model of protein structure such that its lowest normal mode moves a selected set of residues in a targeted way, while still generating realistic and novel proteins. While the current study represents a proof-of-concept with a simple denoiser, relatively short proteins and omitting side-chains, future work is underway to extend this approach to full protein design (Watson et al., 2022; Ingraham et al., 2022) and to investigate whether the so-conditioned samples behave as expected in molecular dynamics simulations. To the best of our knowledge, this is the first time a diffusion model for protein structure has been conditioned on protein dynamics information. We believe this approach holds potential for protein design, where it is desirable to design proteins that are stable and have a desired conformational flexibility.

220 References

- 221 Adamovic, I., Mijailovich, S. M., and Karplus, M. The
222 elastic properties of the structurally characterized
223 myosin ii s2 subdomain: A molecular dynamics
224 and normal mode analysis. *Biophysical Journal*,
225 94(10):3779–3789, 2008. ISSN 0006-3495.
226 doi: <https://doi.org/10.1529/biophysj.107.122028>.
227 URL <https://www.sciencedirect.com/science/article/pii/S0006349508703813>.
- 230 Bahar, I., Atilgan, A. R., and Erman, B. Direct evaluation of
231 thermal fluctuations in proteins using a single-parameter
232 harmonic potential. *Folding and Design*, 2(3):173–181,
233 1997.
- 234 Bahar, I., Lezon, T. R., Bakan, A., and Shrivastava, I. H.
235 Normal Mode Analysis of Biomolecular Structures: Func-
236 tional Mechanisms of Membrane Proteins. *Chemical
237 Reviews*, 110(3):1463–1497, 2010.
- 238 Dhariwal, P. and Nichol, A. Diffusion models beat gans on
239 image synthesis. *ArXiv*, abs/2105.05233, 2021.
- 240 Gibrat, J.-F. and Gō, N. Normal mode analysis of human
241 lysozyme: study of the relative motion of the two domains
242 and characterization of the harmonic motion. *Proteins:
243 Structure, Function, and Bioinformatics*, 8(3):258–279,
244 1990.
- 245 Hinsen, K. and Kneller, G. R. A simplified force field for de-
246 scribing vibrational protein dynamics over the whole fre-
247 quency range. *The Journal of Chemical Physics*, 111(24):
248 10766–10769, 12 1999. ISSN 0021-9606. doi: 10.1063/
249 1.480441. URL <https://doi.org/10.1063/1.480441>.
- 250 Ho, J., Jain, A., and Abbeel, P. Denoising Diffusion Probabilistic Models. *arXiv*, 2020.
- 251 Hoogeboom, E., Satorras, V. G., Vignac, C., and Welling,
252 M. Equivariant diffusion for molecule generation
253 in 3D. In Chaudhuri, K., Jegelka, S., Song, L.,
254 Szepesvari, C., Niu, G., and Sabato, S. (eds.), *Proceedings
255 of the 39th International Conference on Machine Learning*,
256 volume 162 of *Proceedings of Machine Learning Research*, pp. 8867–8887. PMLR, 17–23 Jul
257 2022. URL <https://proceedings.mlr.press/v162/hoogeboom22a.html>.
- 258 Ingraham, J., Baranov, M., Costello, Z., Frappier, V., Ismail,
259 A., Tie, S., Wang, W., Xue, V., Obermeyer, F., Beam,
260 A., and Grigoryan, G. Illuminating protein space with a
261 programmable generative model. *biorXiv*, 2022.
- 262 Jing, B., Eismann, S., Soni, P. N., and Dror, R. O. Equivari-
263 ant graph neural networks for 3d macromolecular struc-
264 ture, 2021.
- 265 Kong, Z., Ping, W., Huang, J., Zhao, K., and Catanzaro, B.
266 Diffwave: A versatile diffusion model for audio synthesis.
267 In *International Conference on Learning Representations*,
268 2021. URL <https://openreview.net/forum?id=a-xFK8Ymz5J>.
- 269 Labesse, G., Colloc'h, N., Pothier, J., and Mornon, J.-P.
270 P-sea: a new efficient assignment of secondary structure
271 from α trace of proteins. *Bioinformatics*, 13(3):291–295,
272 1997.
- 273 Levitt, M., Sander, C., and Stern, P. S. Protein
274 normal-mode dynamics: Trypsin inhibitor, crambin,
275 ribonuclease and lysozyme. *Journal of Molecular
276 Biology*, 181(3):423–447, 1985. ISSN 0022-2836.
277 doi: [https://doi.org/10.1016/0022-2836\(85\)90230-X](https://doi.org/10.1016/0022-2836(85)90230-X).
278 URL <https://www.sciencedirect.com/science/article/pii/002228368590230X>.
- 279 Orengo, C. A., Michie, A. D., Jones, S., Jones, D. T.,
280 Swindells, M. B., and Thornton, J. M. Cath—a hierar-
281 chic classification of protein domain structures. *Structure*,
282 5(8):1093–1109, 1997.
- 283 Schneuring, A., Du, Y., Harris, C., Jamasb, A. R., Igashov,
284 I., weitao Du, Blundell, T. L., Lio, P., Gomes, C. P.,
285 Welling, M., Bronstein, M. M., and Correia, B. Structure-
286 based drug design with equivariant diffusion models,
287 2023. URL <https://openreview.net/forum?id=uKmuzIuVl8z>.
- 288 Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and
289 Ganguli, S. Deep unsupervised learning using nonequi-
290 librium thermodynamics. In Bach, F. and Blei, D. (eds.),
291 *Proceedings of the 32nd International Conference on Ma-
292 chine Learning*, volume 37 of *Proceedings of Machine
293 Learning Research*, pp. 2256–2265, Lille, France, 07–
294 09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/sohl-dickstein15.html>.
- 295 Song, J., Meng, C., and Ermon, S. Denoising diffusion
296 implicit models. *arXiv:2010.02502*, October 2020. URL
297 <https://arxiv.org/abs/2010.02502>.
- 298 Tama, F. and Sanejouand, Y. Conformational change of
299 proteins arising from normal mode calculations. *Protein
300 engineering*, 14 1:1–6, 2001.
- 301 Tirion. Large amplitude elastic motions in proteins from a
302 single-parameter, atomic analysis. *Physical review letters*,
303 77 9:1905–1908, 1996.
- 304 Voet, D. and Voet, J. G. *Biochemistry*. John Wiley & Sons,
305 2010.
- 306 Watson, J. L., Juergens, D., Bennett, N. R., Trippe, B. L.,
307 Yim, J., Eisenach, H. E., Ahern, W., Borst, A. J., Ragotte,
308 R. J., Milles, L. F., Wicky, B. I. M., Hanikel, N., Pellock,

275 S. J., Courbet, A., Sheffler, W., Wang, J., Venkatesh, P.,
276 Sappington, I., Torres, S. V., Lauko, A., Bortoli, V. D.,
277 Mathieu, E., Barzilay, R., Jaakkola, T. S., DiMaio, F.,
278 Baek, M., and Baker, D. Broadly applicable and accu-
279 rate protein design by integrating structure prediction
280 networks and diffusion generative models. *biorXiv*, 2022.

281
282 Zhang, Y. and Skolnick, J. Tm-align: a protein structure
283 alignment algorithm based on the tm-score. *Nucleic acids*
284 *research*, 33(7):2302–2309, 2005.

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330 A. Why Normal Mode Analysis (NMA)?

331 NMA offers unique advantages for generative AI-driven protein design. First, it presents a fast route to a rough hypothesis of
332 *dynamics* for any protein purely based on structure, without the need for expensive molecular dynamics information (Bahar
333 et al., 2010). Second, despite the fact that NMA is less sophisticated than molecular dynamics simulations, the lowest (5-15)
334 non-trivial normal modes have been shown to capture a substantial portion of functional motions in proteins (Bahar et al.,
335 2010). Prominent examples include dihydrofolate reductase (DHFR) (Bahar et al., 1997), lysozyme (Gibrat & Gō, 1990),
336 and adenylate kinase (AdK) (Tama & Sanejouand, 2001), among others.

338 B. What are NMAs assumptions, limitations and merits?

339 The simplicity of NMA comes from its strong assumptions (Bahar et al., 2010). It assumes that the given protein structure
340 represents a potential energy minimum, is in thermal equilibrium, and ignores solvent effects. Despite its simplicity, NMA
341 has been a key tool in protein dynamics research since the 1970s (Bahar et al., 2010). Its successful application spans a
342 variety of proteins, providing key insights into their functional mechanisms. For instance, it elucidated the hinge-bending
343 motion in trypsin (Levitt et al., 1985), and the large-scale conformational changes in myosin (Adamovic et al., 2008).
344 Moreover, NMA has been instrumental in decoding the opening and closing mechanism of AdK (Tama & Sanejouand,
345 2001), and the functional dynamics of DHFR (Bahar et al., 1997). Empirical studies such as these consistently found that the
346 lowest few (typically 5-15) normal modes often capture the predominant part of the functional motions in proteins, beyond
347 the theoretical confines of the harmonic approximation. This makes NMA a practical approach to distill complex protein
348 dynamics into a manageable set of normal modes.

349 There is an array of NMA models, from coarse-grained, where residues or domains are nodes (Tirion, 1996), to highly
350 detailed models, with atoms as nodes and structures refined in molecular force fields (Bahar et al., 2010). Despite the diverse
351 methods, a common finding is that the specific choice of the NMA model does not radically influence the general trend about
352 protein dynamics (Bahar et al., 2010). This suggests that the essential dynamics of proteins can be captured with reasonable
353 accuracy regardless of the model choice such that our method can be adapted in problems requiring complex force-fields.

354 We refer the reader to (Bahar et al., 2010) for a comprehensive review of NMA and its applications.

355 C. Training data

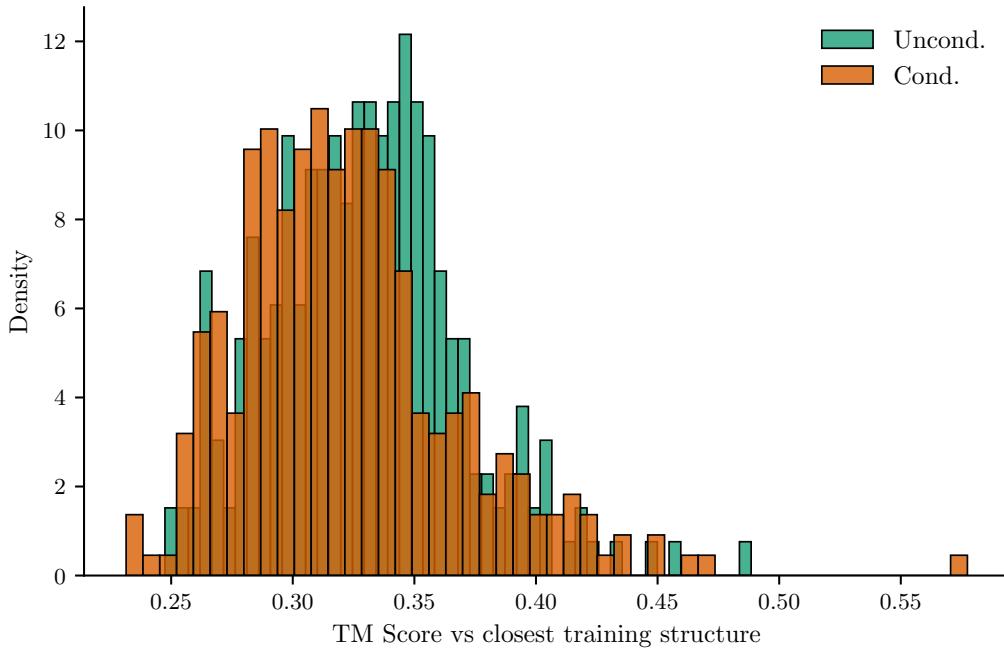
356 Since this work is meant of a proof of concept, we restricted ourselves to short protein sequences to work with limited
357 computational resources. To obtain an interesting sample of short protein snippets, we filtered CATHv4.3 domains (Orengo
358 et al., 1997) for structures with high resolution ($< 3\text{\AA}$), between 20-100 amino acids long. To remove redundancy, we
359 clustered the sequences at 95% sequence similarity. The resulting dataset contains 9'220 protein structures. For these, we
360 set aside 300 for extracting normal mode condition targets that could not have been observed in the training data and use the
361 remaining 8'920 for training.

362 D. Pre-filtering of samples

363 We then generate 3 unconditioned and 4 NMA-conditioned samples for each target condition. We then filter the samples for
364 realistic backbones, requiring that the mean C_α distance along the backbone is within 0.05 \AA of the average C_α distance
365 of 3.8 \AA (Voet & Voet, 2010). This reduces the number of samples by approximately 25% and 50% for unconditional
366 and conditional samples respectively. Of these, we select the sample with the lowest NMA loss (Eq. 4), leaving with
367 approximately 300 samples per sampling procedure.

368 E. Novelty of samples

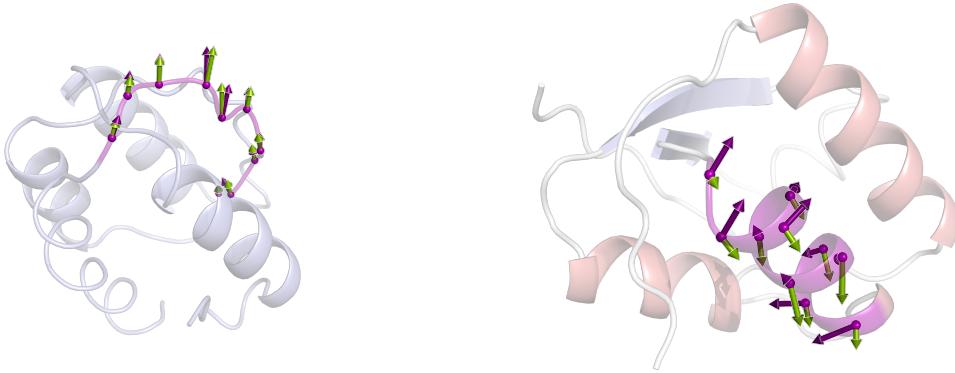
369 To ensure that the model did not overfit we compute the TM-score (Zhang & Skolnick, 2005) between the generated samples
370 and the best matching target fragment of the same size in the training set (Fig. 6, App.). Both, conditioned and unconditioned
371 samples, have TM-scores of about 0.3-0.4, indicating that they are novel compared to the training set and do not simply
372 memorise the training data. Importantly, NMA-conditioning does not have a negative effect on the novelty of the samples.



407 Figure 6: Histogram of the TM-score between the generated samples and the best matching target in the training set. With
408 TM scores around 0.3-0.4, unconditioned and conditioned samples both generate novel samples and NMA-conditioning
409 does not have a negative effect on novelty.

410 411 F. Extra samples 412

413 Below we provide a few more samples from the unconditional and conditional sampling process, to give a better idea of the
414 diversity and quality of the samples.
415



429 Figure 7: **Left:** conditioned sample with L_1 loss 0.05. **Right:** not conditioned sample L_1 loss 0.89.
430
431
432
433
434
435
436
437
438
439

440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459

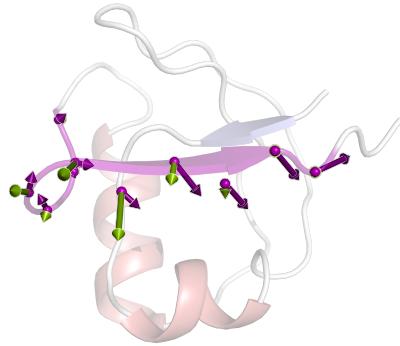
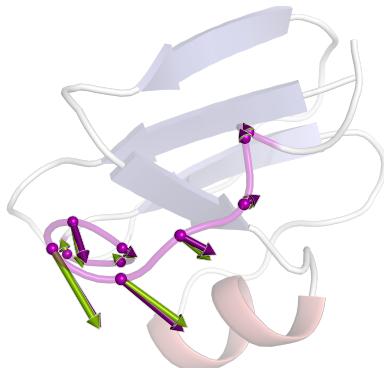


Figure 8: **Left:** conditioned sample with L_1 loss 0.04. **Right:** not conditioned sample L_1 loss 1.02

460
461
462
463
464
465
466
467
468
469
470
471
472
473

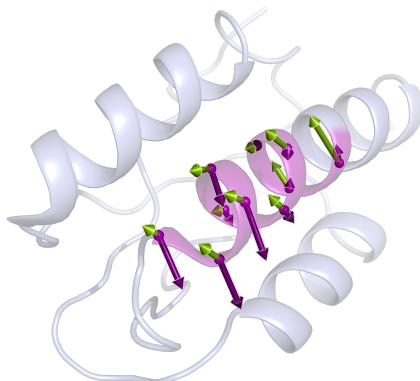
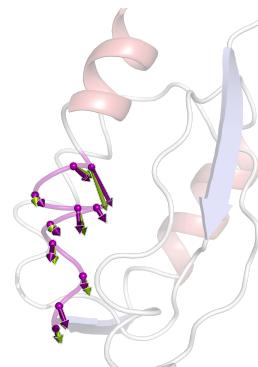


Figure 9: **Left:** conditioned sample, loss 010. **Right:** not conditioned sample, loss 0.82

487
488
489
490
491
492
493
494