

1 **Title**

2 Cyclic peptides space: The methodology of sequence selection to cover the comprehensive physical  
3 properties

4  
5 **Authors**

6 Ryo Tsuchihashi<sup>1\*</sup>, Misaki Kinoshita<sup>1\*</sup> (These authors contributed equally to this work.)

7 1. Research & Development Division, Central Research Laboratory, Japan Blood Products Organization,  
8 Kobe, Japan

9 (\*Corresponding authors: Ryo Tsuchihashi, [tsuchihashi-ryo@jbpo.or.jp](mailto:tsuchihashi-ryo@jbpo.or.jp), Misaki Kinoshita, [kinoshita-](mailto:kinoshita-misaki@jbpo.or.jp)  
10 [misaki@jbpo.or.jp](mailto:misaki@jbpo.or.jp))

11  
12 **Abstract**

13 Cyclic peptides have emerged as a pivotal modality for next-generation therapeutics, due to their superior  
14 biocompatibility, high selectivity, and structural stability. While AI-driven peptide design has advanced rapidly,  
15 conventional optimization algorithms are often constrained by initialization biases, which impede the efficient  
16 exploration of the vast chemical space. Here, we propose a novel methodology that integrates the protein language  
17 model ESM-2 with cyclic permutation averaging of embeddings to resolve this bottleneck. This approach  
18 establishes a comprehensive "peptide space", a high-dimensional vector representation that encapsulates the  
19 physicochemical and structural attributes of cyclic peptides. Our analysis reveals that random sequence selection  
20 results in a heterogeneous distribution within this space, potentially underrepresenting specific functional regions.  
21 Conversely, navigating this defined peptide space enables the selection of libraries that uniformly span diverse  
22 molecular properties. In a proof-of-concept study designing binders for  $\beta$ 2-microglobulin ( $\beta$ 2m), we demonstrate  
23 that initial sequences uniformly sampled from our peptide space yield superior candidates more efficiently than  
24 those derived from random selection. Furthermore, this framework facilitates the quantitative assessment of  
25 mutational perturbations on global peptide properties, supporting rational decision-making for both broad  
26 exploration and local optimization. This "peptide space" concept provides a foundational framework for defining  
27 appropriate search boundaries and enhancing computational efficiency in AI-mediated drug discovery.

28  
29 **Keywords**

30 Cyclic peptide, ESM-2, Peptide design, Drug discovery

1

2 **Acknowledgements**

3 -

4

1

## 2 **Introduction**

3           In the landscape of pharmaceutical development, peptide therapeutics have emerged as a pivotal modality,  
4 bridging the gap between conventional small molecules and antibody biologics [1, 2]. Peptides offer a unique  
5 advantage by accessing complex chemical spaces and adopting intricate conformations, thereby enabling high-  
6 affinity interactions with a broad spectrum of disease-associated proteins which are often considered  
7 "undruggable" by small molecules. Furthermore, they can be engineered for specific properties such as enhanced  
8 tissue penetration. Compared to antibodies, peptides possess superior biocompatibility due to their construction  
9 from fundamental amino acids and offer the distinct benefit of cost-effective synthetic manufacturability [3,  
10 4]. Consequently, the market for peptide therapeutics is expanding rapidly, with successive approvals in diverse  
11 therapeutic areas including diabetes, oncology, cardiovascular diseases, and rare disorders [5]. Among these,  
12 cyclic peptides represent a privileged scaffold. Beyond the general advantages of ease of synthesis and high target  
13 selectivity, macrocyclization confers critical physicochemical benefits: resistance to peptidases (proteolytic  
14 stability), reduced conformational entropy leading to higher binding affinity, and improved membrane  
15 permeability [3]. While recent advancements in AI-driven design have significantly enhanced the accuracy and  
16 productivity of cyclic peptide engineering, the comprehensive exploration of the vast combinatorial search space  
17 remains computationally prohibitive [6]. To mitigate this, existing computational frameworks employ heuristic  
18 optimization algorithms to accelerate the search process. For instance, recent systems such as HighPlay [7] and  
19 EvoBind2 [8] utilize evolutionary algorithms initiated from arbitrarily determined sequences. However, because  
20 the initial seed sequences profoundly influence the trajectory and quality of the final solution in these stochastic  
21 methods, defining an appropriate "search space" for initialization is crucial. Yet, systematic approaches to define  
22 such spaces have been lacking. Although some approaches, like RFpeptide, attempt to address this by utilizing a  
23 predefined initial space of 1,200 sequences rather than a single seed [9], this space is constructed with a primary  
24 focus on secondary structure characteristics, largely overlooking the diversity of chemical properties.

25           In this study, we report the construction of a comprehensive design space that facilitates the rational  
26 determination of initial sequences, incorporating a broader range of physicochemical and structural attributes. We  
27 generated a high-dimensional vector space by transforming random amino acid sequences using the protein  
28 language model ESM-2 [10]. We evaluated this space for its uniformity in physicochemical properties and its  
29 utility in binder design. Our results demonstrate that this peptide space enables the unbiased selection of properties  
30 and serves as a powerful tool for optimizing both initial sequence selection and evolutionary directionality. This

1 framework provides a novel methodology for the efficient and appropriate selection of cyclic peptide sequences  
2 through the explicit definition and understanding of the search space.

3

4

## 5 **Methods**

### 6 **Protein Language Model and Embedding Generation**

7 To extract numerical feature vectors (embeddings) from peptide sequences, we employed the pre-trained  
8 protein language model ESM-2 (esm.pretrained.esm2\_t6\_8M\_UR50D) [10]. Specifically, we utilized residue-  
9 level representation vectors derived from the intermediate layer (layer 6) [11]. To ensure deterministic and  
10 reproducible embeddings, the model was operated in evaluation mode, thereby disabling dropout layers.

### 11 **Representation Vectors for Cyclic Peptides**

12 Cyclic peptides lack a defined N- or C-terminus. To incorporate this topological characteristic into the  
13 vector representation, we introduced a "Cyclic Permutation Averaging" strategy. For a peptide sequence  $S$  of  
14 length  $L$ , we first generated all  $L$  possible cyclic permutations (denoted as sequence set  $S_i$ , where  $i = 0, \dots, L - 1$ ,  
15 and  $S_0$  is the original sequence), effectively shifting the sequence one residue at a time:

$$16 \quad S_i = S_0[i:] + S_0[:i]$$

17 Next, for all  $L$  sequences, we computed the representation vector  $R_i$  for the entire sequence using the ESM-2  
18 model described above. We calculated the arithmetic mean of these  $L$  sequence representation vectors to obtain  
19 the topology-invariant embedding vector,  $R_{cyclic}$ , for the original cyclic peptide  $S$ :

$$20 \quad R_{cyclic} = \frac{1}{L} \sum_{i=0}^{L-1} R_i \quad (\text{function 1})$$

21 , where  $R_i$  indicated the vector for the sequence  $S_i$ . As a control, we also defined a conventional linear peptide  
22 representation derived from a single sequence, without applying cyclic permutation averaging.

### 23 **Construction of Peptide Space and Dimensionality Reduction**

24 In this study, we focused on peptides with a length of 14 amino acids. We constructed a large-scale  
25 random peptide library consisting of approximately one million ( $N \sim 300,000$ ) sequences, generated by randomly  
26 sampling from the standard 20 amino acids. High-dimensional embedding vectors were calculated for all  
27 sequences using the Cyclic Permutation Averaging method defined above. To visualize and analyze the global  
28 structure of this high-dimensional embedding landscape, we applied Uniform Manifold Approximation and

1 Projection (UMAP), a non-linear dimensionality reduction technique [12]. UMAP was used to project the high-  
2 dimensional vectors onto a two-dimensional plane. In this manuscript, we refer to this 2D projection as the  
3 "Peptide Space."

#### 4 **Characterization of the Peptide Space**

5 To characterize the landscape of the constructed chemical space, we performed the following three analyses. First,  
6 we quantified the distribution density of peptide sequences within the peptide space using Kernel Density  
7 Estimation (KDE) to visualize potential sampling biases. For bandwidth selection, we employed Scott's Rule,  
8 which automatically optimizes the parameter based on sample size and dimensionality to achieve a balance  
9 between preventing over-smoothing and avoiding noise over-fitting. Second, The abundance of specific amino  
10 acid residues was mapped onto the peptide space. This allowed us to evaluate the correlation between intrinsic  
11 chemical properties and their relative positioning within the projected manifold. Finally, to assess the local and  
12 global organization of the space, we generated two derivative sets from reference sequences: "shuffled sequences"  
13 (preserving composition but altering order) and "variant libraries" (comprising comprehensive residue  
14 substitutions). We then analyzed their spatial distribution and calculated the cosine similarity within the cyclic  
15 peptide space to define the relationship between sequence homology and spatial proximity.

16

#### 17 **Empirical Study: Optimization and Sampling Strategy for $\beta$ 2m binders**

18 To demonstrate the utility of our approach, we analyzed datasets obtained from computational  
19 optimization simulations of binder peptides against  $\beta$ 2m. To evaluate the impact of initial seed selection on  
20 discovery efficiency, we compared two distinct sequence set. One is the Systematic sampling set, in which the  
21 peptide space was partitioned into a uniform grid. We identified 92 specific grids that met a minimum sequence  
22 density threshold. Representative sequences were then uniformly extracted from each of these grids to ensure  
23 maximal coverage of the manifold. Another is the Random sampling set. In this set, sequences were extracted  
24 from the library via stochastic selection, without regard to their spatial distribution or chemical diversity. To  
25 ensure a rigorous comparison, the total number of sequences for both sequence set was fixed at 920, corresponding  
26 to 10 sequences per grid. Using these seed sequences, we conducted binding discovery and structure prediction  
27 simulations employing EvoBind2. The primary metric was the Loss value, a composite score reflecting predicted  
28 binding free energy and structural stability (Equation 2), where lower values indicate higher predicted affinity.  
29 Simulations were organized into 10 independent sets, each comprising 92 sequences (corresponding to the number

1 of grids). The minimum Loss value from each set was extracted to generate box plots for comparative statistical  
2 analysis.

$$3 \quad \text{Loss} = \text{peptide } p\text{LDDT}^{-1} \cdot \left( \frac{1}{n} \sum_{j=1}^n d_j \right) \quad (\text{function 2})$$

#### 4 **Correlation between Peptide Space and Structure/Properties**

5 We evaluated the correlation between coordinates in the peptide space and the structural features and  
6 physicochemical properties of the peptides. Structural predictions were performed using AfCycDesign [15] for  
7 3,000 sequences randomly selected within the peptide space. From the predicted structures, we calculated  
8 topological features—such as radius and inter-atomic distances—secondary structure content, and known  
9 physicochemical properties, including hydrophobicity and charge.

10

### 11 **Results and Discussion**

#### 12 **Construction of a cyclic peptide space for unbiased property selection**

13 Establishing an appropriate initialization strategy for cyclic peptide design requires the construction of a  
14 search space that is sufficiently discrete and representative of diverse physicochemical and structural properties.  
15 We utilized the protein language model ESM-2 [10] to transform randomly generated sequences into high-  
16 dimensional vectors that encapsulate their inherent characteristics. Given that ESM-2 is trained on linear  
17 sequences, we addressed the cyclic nature of our targets by generating cyclic permutations of the amino acid  
18 sequence (shifting the starting position sequentially), computing embeddings for each permutation, and averaging  
19 them to derive a definitive "cyclic peptide vector" (Fig. 1A). We constructed vector spaces for datasets comprising  
20 1,000, 10,000, and 300,000 random 14-residue sequences. Dimensionality reduction via UMAP [12] revealed that  
21 spaces exceeding 10,000 sequences exhibit a characteristic distribution partitioned into three distinct segments.  
22 Notably, random sampling resulted in significant heterogeneity in occurrence frequency within this peptide space  
23 (Fig. 1A, right panels). While similar segmentation is observed in linear sequences (Supporting Figure 1), the  
24 cyclic vectors derived from the 300,000-sequence dataset eliminated several minor clusters observed in linear  
25 counterparts, suggesting that our cyclic permutation averaging effectively smooths out subtle variances caused by  
26 linear edge effects.

27 To validate the robustness of the proposed cyclic permutation averaging method, we evaluated the  
28 variance and cosine similarity of four vector sets derived from a single seed sequence (DTIAVEHDLGAVQE):

1 (i) linear vectors from cyclic permutations; (ii) cyclic vectors from cyclic permutations; (iii) linear vectors from  
2 shuffled sequences; and (iv) cyclic vectors from shuffled sequences (Supporting Figure 2). For the linear vectors  
3 without averaging (i), slight fluctuations were observed depending on the starting amino acid position, resulting  
4 in a cosine similarity of 0.997764 (Supporting Figure 2B, C). This bias is attributed to edge effects (N/C-terminal  
5 recognition) inherent in linear models. In contrast, the vectors obtained through the cyclic permutation averaging  
6 method (ii), converged to nearly identical coordinates in the peptide space regardless of the input permutation,  
7 achieving a perfect cosine similarity of 1.0000. Although these vectors occasionally appeared slightly separated  
8 into two points on the UMAP plot due to computational noise or approximation, they were confirmed to be  
9 effectively identical. Conversely, the shuffled sequences (iii, iv)—which maintain the same amino acid  
10 composition but randomized order—showed significant dispersion in the space for both linear and cyclic formats  
11 (Supporting Figure 2B, C). These results demonstrate that the vectors generated by our method reflect structural  
12 and chemical properties based on sequence order rather than simple composition, thereby appropriately  
13 representing the topology of cyclic peptides.

14 Mapping randomly picked sequences onto this space revealed no correlation between spatial distance  
15 and sequence similarity, indicating that the space captures properties independent of simple sequence homology  
16 (Fig. 1B). Furthermore, to assess the "randomness" relative to a specific query, we calculated cosine similarities  
17 between the seed sequence (DTIAVEHDLGAVQE) and 10,000 random sequences (Fig. 2B). The majority of  
18 random sequences clustered within a high similarity range (0.85–1.00), demonstrating that naive random sampling  
19 frequently results in vectors with redundant properties. Color-coding the distribution by amino acid count revealed  
20 that the three major segments in the peptide space are largely defined by cysteine content (Supporting Figure 3A).  
21 Other residues also drive spatial bias; for instance, peptides constructed with more than two methionine localize  
22 between the central and right segments, while those with constructed with more than two tryptophan show  
23 distinct biases within segments. This underscores that simple amino acid composition significantly skews the  
24 distribution.

25 We further analyzed structural attributes by generating 3D structures for 1,000 randomly selected  
26 sequences using AfCycDesign [15], an AlphaFold-based tool for cyclic peptides. We mapped various structural  
27 or physical features onto the peptide space (Supporting Figure 3B). Rather than a uniform distribution,  
28 these parameters exhibited a mosaic-like, patchy organization. Physicochemical properties showed clear biases,  
29 consistent with the aforementioned compositional skew. Similarly, secondary structure and conformational  
30 metrics displayed non-uniform, localized patterns. Collectively, these findings suggest that random sequence

1 selection fails to guarantee physicochemical or structural stochasticity, potentially making it an unsuitable strategy  
2 for initialization. Our proposed "peptide space" offers a solution to overcome these biases by enabling rational,  
3 uniform sampling.

4

#### 5 **Use case: Enhanced binder design via peptide space navigation**

6 To demonstrate the utility of this peptide space, we applied it to the design of cyclic peptide binders for  
7  $\beta$ 2m using EvoBind2 [8].  $\beta$ 2m is a component of MHC class I molecules, essential for presenting self-antigens to  
8 cytotoxic T cells [16]. With its seven-stranded  $\beta$ -sandwich fold and single disulfide bond within a compact 99-  
9 residue structure,  $\beta$ 2m represents an ideal target for evaluating complex prediction. To ensure uniform property  
10 sampling, we segmented the peptide space into 92 grids and randomly selected 10 sequences from each,  
11 creating a "UMAP-based sequence set" (920 sequences). A control "Random sequence set" of equal size was  
12 generated stochastically (Fig. 2A). The optimization process involved 10 iterations, each using a subset of 92  
13 sequences (one from each grid for the UMAP-based set, and a random subset for the control).

14 The grid-based selection demonstrated a broader distribution of cosine similarities compared to random  
15 selection, specifically increasing the probability of sampling sequences with lower similarity (i.e., unique  
16 properties) (Fig. 2B). We then performed EvoBind2 design targeting  $\beta$ 2m for each subset and calculated the Loss  
17 value using the equation 2 described in the methods section. The UMAP-based sequence set yielded consistently  
18 lower mean and minimum Loss values compared to the Random sequence set, indicating a higher success rate in  
19 identifying computationally favorable candidates (Fig. 2C). Furthermore, mapping the loss values from the  
20 UMAP-based set in the peptide space revealed that the sequences exhibiting the lowest loss values are located at  
21 the segment boundaries (Supplementary Figure 4), because this region is often missed by random sampling due  
22 to its statical rarity. These results demonstrate that navigating this peptide space ensures true physicochemical  
23 coverage, which enabled the efficient identification of high-potential candidates that would otherwise be  
24 overlooked.

25

#### 26 **Quantifying mutational effects on peptide properties**

27 Finally, we utilized this space to quantitatively evaluate how specific amino acid substitutions perturb  
28 global peptide properties. Using the sequence DTIAVEHDLGAVQE, we substituted the residues at positions 2  
29 and 6 with all 20 amino acids and analyzed the resulting shifts in the peptide space and cosine distances. Single  
30 substitutions at either position caused minimal spatial displacement, with the notable exception of cysteine

1 (Fig. 3A, B). Cysteine introduction triggered a segment shift in peptide space, reflecting the dominant influence  
2 of disulfide potential (Supporting Figure 3A). While non-cysteine mutations showed no consistent directional  
3 trend in peptide space (Fig. 3A, B, right panels), hierarchical clustering based on cosine similarity revealed clear  
4 physicochemical logic (Supporting Figure 5). This discrepancy arises because UMAP dimensionality reduction  
5 preserves local topology but may distort global directional relationships. Clustering analysis confirmed that  
6 substitutions with chemically similar residues—such as Asp/Glu, Ser/Thr, Tyr/Phe/Ile, and Ala/Val—formed tight  
7 clusters with high similarity. This is intuitive given their side-chain properties. Interestingly, basic residues (Lys,  
8 Arg, His) did not form a unified cluster, indicating that in the context of cyclic peptides, these residues impart  
9 distinct vector characteristics, likely due to differences in functional group geometry and hydrophobic surface  
10 area. Double mutations (positions 2 and 6) resulted in a synergistic expansion of distribution in peptide space and  
11 high variance in cosine similarity (Fig.3C, D). This highlights two key capabilities: (1) the ability to switch  
12 segments by altering cysteine content, and (2) the ability to fine-tune similarity within a segment by combining  
13 specific mutations. In conclusion, this peptide space provides a powerful framework for rational design, enabling  
14 users to select mutations that match their specific exploration goals—whether it be a broad "jump" to a new  
15 physicochemical regime or a local "optimization" preserving the parent character.

16

## 17 **Conclusion**

18 In this study, we proposed the novel concept of "Peptide Space." This framework employs a strategy  
19 where linear sequence embeddings generated by ESM-2 are averaged over a set of cyclic permutations, explicitly  
20 designed to mitigate the bias of terminal effects inherent in language models trained on linear proteins. Our  
21 analysis demonstrated that naive random selection of sequences results in a non-uniform distribution of  
22 physicochemical and structural properties within the ESM-2 space. This finding underscores a critical discrepancy  
23 in conventional machine learning-based design: while the objective of initialization is to cover a "random" spread  
24 of physicochemical properties, relying on "sequence" randomness leads to the systematic undersampling of low-  
25 frequency, yet potentially functional, regions of the chemical space. Our methodology resolves this bottleneck,  
26 enabling a truly unbiased initialization that encompasses a broad spectrum of molecular attributes.

27 In our proof-of-concept study designing binders for  $\beta 2m$  using EvoBind2, we compared the binding  
28 potential, quantified by Loss values, of candidates derived from our peptide space against those from a  
29 random sequence set. The results indicated that navigating the peptide space allows for the identification of  
30 superior starting points. This observation likely extends beyond  $\beta 2m$  to other therapeutic targets; reliance on

1 stochastic sequence initialization carries a high probability of overlooking sequences with favorable theoretical  
2 binding probability (low Loss values), thereby compromising the quality of the final design and inflating  
3 computational costs. Conversely, utilizing our peptide space facilitates the selection of appropriate initial seeds,  
4 helping to circumvent local minima, enhance the quality of the final design, and reduce the overall computational  
5 burden.

6 From the perspective of computational efficiency, this framework is particularly advantageous for  
7 modulating evolutionary directionality. By utilizing changes in the direction and distance of peptide sequence  
8 vectors within reduced-dimension coordinates, it becomes possible to rationally select next-generation  
9 mutations tailored to specific local optimization goals. For instance, in purpose of excluding mutations that result  
10 in negligible vector shifts involving chemically minor changes, efficiency can be significantly improved by  
11 pruning the search space rather than relying on blind stochastic mutation. Although integrating ESM-2  
12 calculations adds a step to the peptide design workflow, the inference time is minimal (<1 ms per sequence [10])  
13 and is far outweighed by the efficiency gains achieved through the strategic reduction of the search space.

14 A fundamental insight of this work is the critical importance of comprehending the nature of the search  
15 space itself. As the application of machine learning expands from cyclic peptides to broader domains of drug  
16 discovery and materials science—including the design of protein binders and inhibitors [17]—the associated  
17 computational costs are becoming a non-negligible barrier. Our approach, while conceptually simple,  
18 demonstrates that explicitly defining and understanding the search space can lead to substantial cost reductions  
19 and more rapid, robust design cycles. When the concept of "Chemical Space" was first introduced, its primary  
20 utility was the comprehensive coverage of features [18, 17]. However, in the current era of generative AI, we  
21 argue that this concept has evolved to play a dual role: ensuring diversity and optimizing computational tractability.  
22 As new machine learning architectures for peptides and proteins continue to emerge, methods that provide a  
23 structured understanding of the exploration landscape—such as the one proposed here—will play an indirect but  
24 indispensable role in maximizing the efficiency of our limited computational resources.

25

26

## 27 **Statements and Declarations**

### 28 **Funding**

29 The authors declare that no funds, grants, or other support were received during the preparation of this manuscript.

30

## 1 **Competing Interests**

2 The authors have no relevant financial or non-financial interests to disclose.

3

## 4 **Ethics Approval**

5 This article does not contain any studies with human participants or animals performed by any of the authors.

6

## 7 **Data Availability**

8 All data generated or analyzed during this study are included in this published article and its supplementary  
9 information files.

10

## 11 **References**

12 [1] Xiao W, Jiang W, Chen Z, Huang Y, Mao J, Zheng W, Hu Y, Shi J (2025) Advance in peptide-based drug  
13 development: delivery platforms, therapeutics and vaccines. *Sig Transduct Target Ther* 10:.  
14 <https://doi.org/10.1038/s41392-024-02107-5>

15 [2] Rossino G, Marchese E, Galli G, Verde F, Finizio M, Serra M, Linciano P, Collina S (2023) Peptides as  
16 Therapeutic Agents: Challenges and Opportunities in the Green Transition Era. *Molecules* 28:7165.  
17 <https://doi.org/10.3390/molecules28207165>

18 [3] Ji X, Nielsen AL, Heinis C (2024) Cyclic Peptides for Drug Development. *Angew Chem Int Ed* 63:.  
19 <https://doi.org/10.1002/anie.202308251>

20 [4] Lamers C (2022) Overcoming the Shortcomings of Peptide-Based Therapeutics. *Future Drug. Discov.* 4:.  
21 <https://doi.org/10.4155/fdd-2022-0005>

22 [5] de la Torre BG, Albericio F (2020) Peptide Therapeutics 2.0. *Molecules* 25:2293.  
23 <https://doi.org/10.3390/molecules25102293>

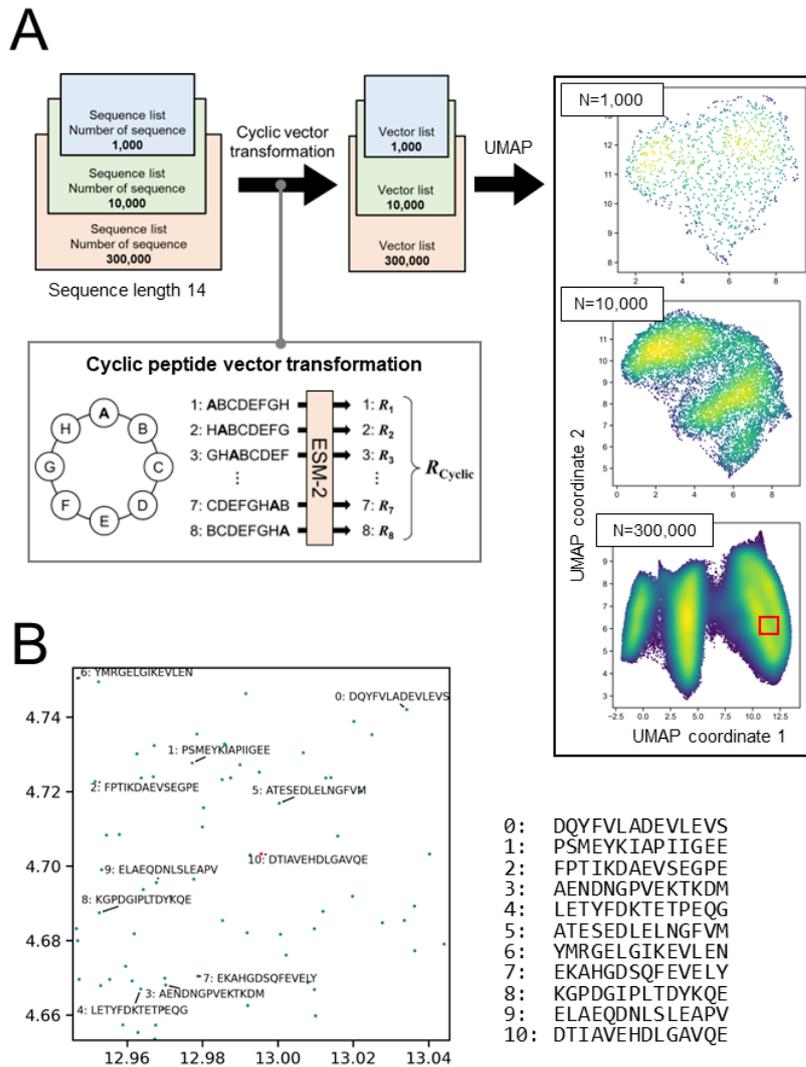
24 [6] Joo S (2012) Cyclic Peptides as Therapeutic Agents and Biochemical Tools. *Biomolecules and Therapeutics*  
25 20:19-26. <https://doi.org/10.4062/biomolther.2012.20.1.019>

26 [7] Lin H, Zhu C, Shang T, Zhu N, Lin K, Zhang C, Shao X, Wang X, Duan H (2025) HighPlay: Cyclic Peptide  
27 Sequence Design Based on Reinforcement Learning and Protein Structure Prediction. *J. Med. Chem.* 68:12047-  
28 12057. <https://doi.org/10.1021/acs.jmedchem.5c00896>

29 [8] Li Q, Vlachos EN, Bryant P (2025) Design of linear and cyclic peptide binders from protein sequence  
30 information. *Commun Chem* 8:.  
<https://doi.org/10.1038/s42004-025-01601-3>

- 1 [9] Rettie SA, Juergens D, Adebomi V, Bueso YF, Zhao Q, Leveille AN, Liu A, Bera AK, Wilms JA, Üffing A,  
2 Kang A, Brackenbrough E, Lamb M, Gerben SR, Murray A, Levine PM, Schneider M, Vasireddy V, Ovchinnikov  
3 S, Weiergräber OH, Willbold D, Kritzer JA, Mougous JD, Baker D, DiMaio F, Bhardwaj G (2025) Accurate de  
4 novo design of high-affinity protein-binding macrocycles using deep learning. *Nat Chem Biol* 21:1948-1956.  
5 <https://doi.org/10.1038/s41589-025-01929-w>
- 6 [10] Lin Z, Akin H, Rao R, Hie B, Zhu Z, Lu W, Smetanin N, Verkuil R, Kabeli O, Shmueli Y, Dos Santos Costa  
7 A, Fazel-Zarandi M, Sercu T, Candido S, Rives A (2023) ,Evolutionary-scale prediction of atomic-level protein  
8 structure with a language model.*Science*379,1123-1130. <https://doi.org/10.1126/science.ade2574>
- 9 [11] Frank M, Ni P, Jensen M, Gerstein MB (2024) Leveraging a large language model to predict protein phase  
10 transition: A physical, multiscale, and interpretable approach. *Proc. Natl. Acad. Sci. U.S.A.* 121:.  
11 <https://doi.org/10.1073/pnas.2320510121>
- 12 [12] McInnes L, Healy J, Melville J (2020) UMAP: Uniform Manifold Approximation and Projection for  
13 Dimension Reduction. *arXiv:1802.03426v3*. <https://doi.org/10.48550/arXiv.1802.03426>
- 14 [13] Mariani V, Biasini M, Barbato A, Schwede T (2013) IDDT: a local superposition-free score for comparing  
15 protein structures and models using distance difference tests. *Bioinformatics* 29:2722-2728.  
16 <https://doi.org/10.1093/bioinformatics/btt473>
- 17 [14] Zhang Y (2005) TM-align: a protein structure alignment algorithm based on the TM-score. *Nucleic Acids*  
18 *Research* 33:2302-2309. <https://doi.org/10.1093/nar/gki524>
- 19 [15] Rettie SA, Campbell KV, Bera AK, Kang A, Kozlov S, Bueso YF, De La Cruz J, Ahlrichs M, Cheng S,  
20 Gerben SR, Lamb M, Murray A, Adebomi V, Zhou G, DiMaio F, Ovchinnikov S, Bhardwaj G (2025) Cyclic  
21 peptide structure prediction and design using AlphaFold2. *Nat Commun* 16:.  
22 <https://doi.org/10.1038/s41467-025-59940-7>
- 23 [16] Wu Y, Zhang N, Hashimoto K, Xia C, Dijkstra JM (2021) Structural Comparison Between MHC Classes I  
24 and II; in Evolution, a Class-II-Like Molecule Probably Came First. *Front. Immunol.* 12:.  
25 <https://doi.org/10.3389/fimmu.2021.621153>.
- 26 [17] Lipinski CA, Lombardo F, Dominy BW, Feeney PJ (1997) Experimental and computational approaches to  
27 estimate solubility and permeability in drug discovery and development settings. *Advanced Drug Delivery*  
28 *Reviews* 23:3-25. [https://doi.org/10.1016/S0169-409X\(96\)00423-1](https://doi.org/10.1016/S0169-409X(96)00423-1)
- 29 [18] Dobson CM (2004) Chemical space and biology. *Nature* 432:824-828. <https://doi.org/10.1038/nature03192>
- 30

1 **Figures**



2

3 **Fig. 1 | Construction and characterization of a latent space for cyclic peptides.** (A) Schematic workflow for

4 generating cyclic peptide embeddings. Randomly generated peptide sequences (length 14) with varying dataset

5 sizes ( $N = 1,000, 10,000,$  and  $300,000$ ) were transformed into high-dimensional vectors using the ESM-2. To

6 accommodate the circular topology, the final representation ( $R_{Cyclic}$ ) was computed by aggregating embeddings

7 derived from all cyclic permutations of the linear sequence (inset). Dimensionality reduction via UMAP (right

8 panels) reveals that a characteristic tripartite distribution emerges as the sample size exceeds 10,000. (B) Local

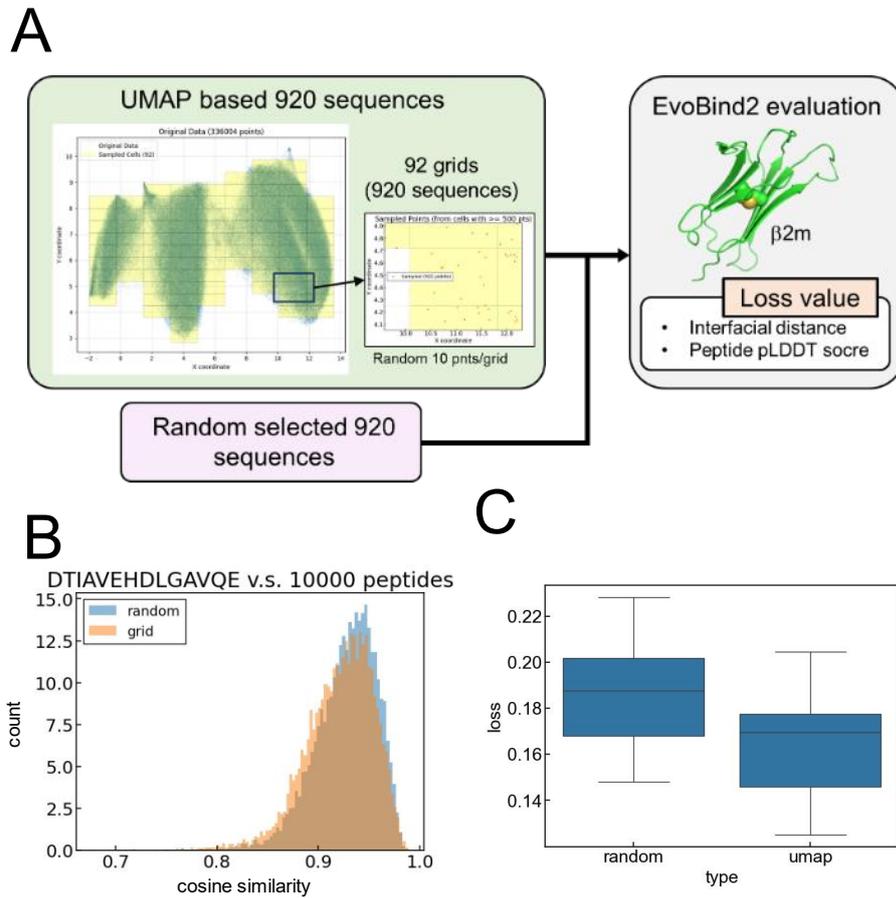
9 analysis of sequence-space relationships. A magnified view of the high-density region (indicated by the red square

10 in A) shows the projection of specific peptide sequences with sequence list (right) corresponding to the scatter

11 points. Since the peptide space is defined by the UMAP-based coordinate system, the x and y axis labels are

12 omitted.

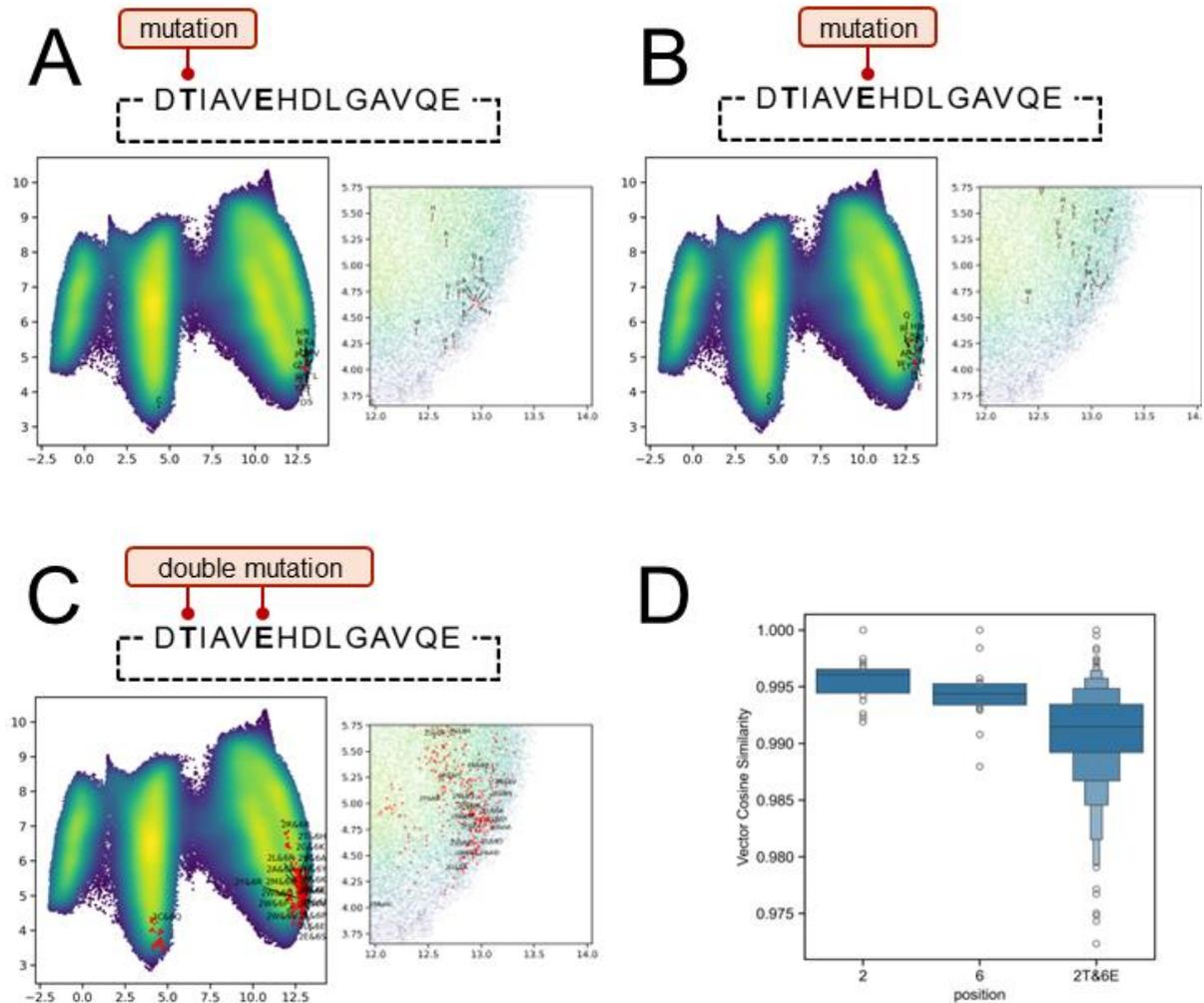
1



2

3 **Fig. 2 | Enhanced efficiency of cyclic peptide design via latent space-guided initialization.** (A) Workflow for  
4 evaluating the utility of the constructed peptide space in targeted design. The target protein,  $\beta 2m$ , was used to test  
5 the generation of cyclic peptide binders via EvoBind2 [8]. To mitigate sampling bias, the peptide space was  
6 segmented into 92 discrete grids, and 10 sequences were randomly drawn from each to create a structurally diverse  
7 "UMAP-based" dataset (N=920). This was compared against a control "Random" dataset (N=920) generated  
8 stochastically. Loss values, calculated based on interfacial distance and peptide pLDDT scores (inset), served as  
9 the evaluation metric. (B) Assessment of sequence diversity. The histogram displays the distribution of cosine  
10 similarities for the UMAP-based (orange) versus random (blue) datasets against a reference sequence. The grid-  
11 based selection exhibits a broader distribution with a higher frequency of lower similarity scores, indicating that  
12 spatial segmentation effectively reduces redundancy and captures a wider range of physicochemical properties.  
13 (C) Comparison of design performance. Box plots representing the distribution of loss values after EvoBind2  
14 optimization.

15



1

2 **Fig. 3 | Visualization and quantification of mutational effects in the latent space.** (A–C) Projection of mutant  
3 sequences onto the peptide space. The reference sequence was subjected to exhaustive point mutations at position  
4 2 (A), position 6 (B), or simultaneous mutations at both positions (C). In each panel, the left plot displays the  
5 global distribution of the resulting mutant vectors within the constructed UMAP landscape, while the right plot  
6 provides a magnified view of the local region surrounding the parental sequence to visualize the relative  
7 positioning of individual variants. (D) Distribution of vector similarities. Box plots show the cosine similarity of  
8 single ("2" and "6") and double ("2T&6E") mutant vectors relative to the reference sequence. Since the peptide  
9 space is defined by the UMAP-based coordinate system, the x and y axis labels are omitted.

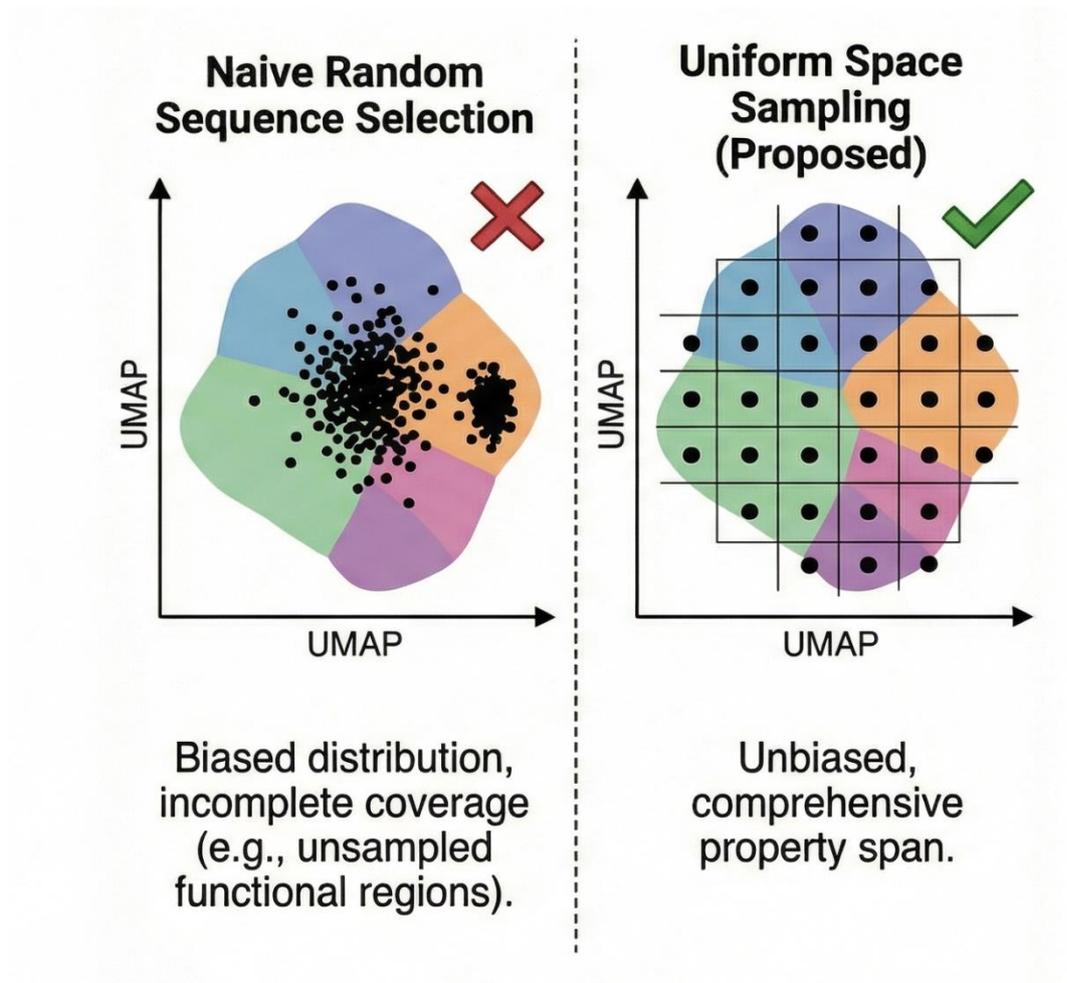
10

11

12

13

1 Graphic Abstract  
2

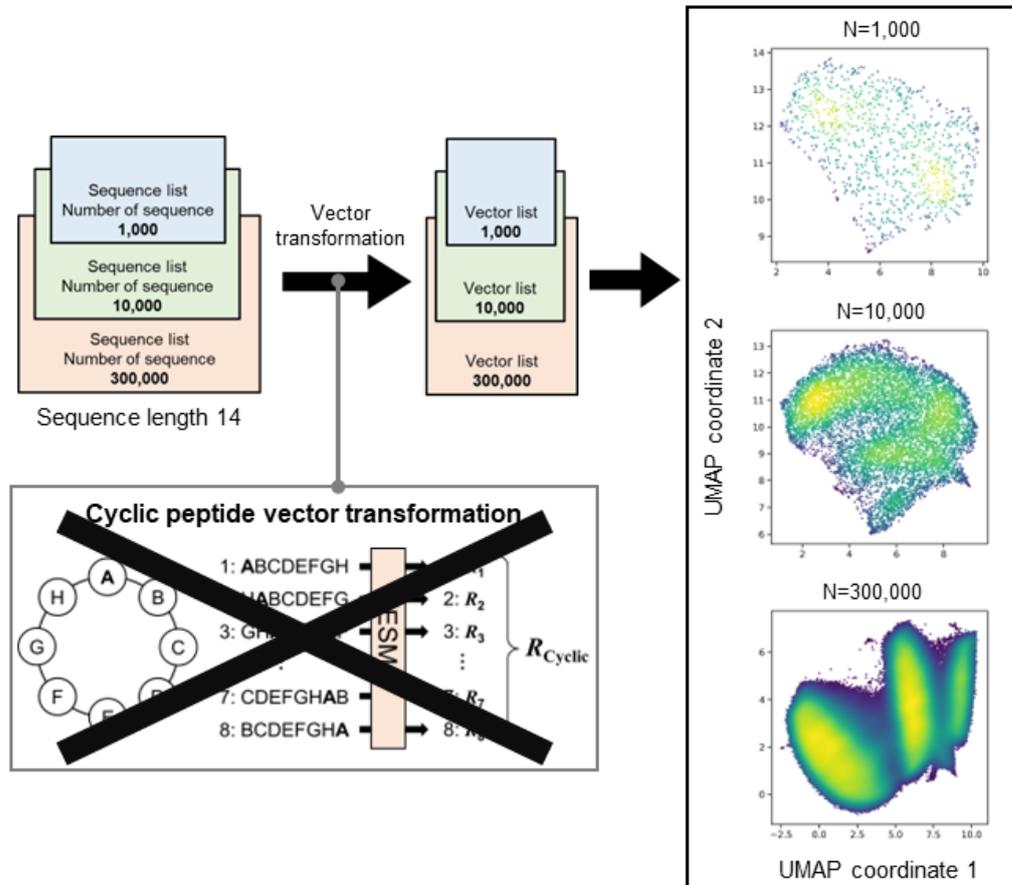


3  
4  
5  
6  
7

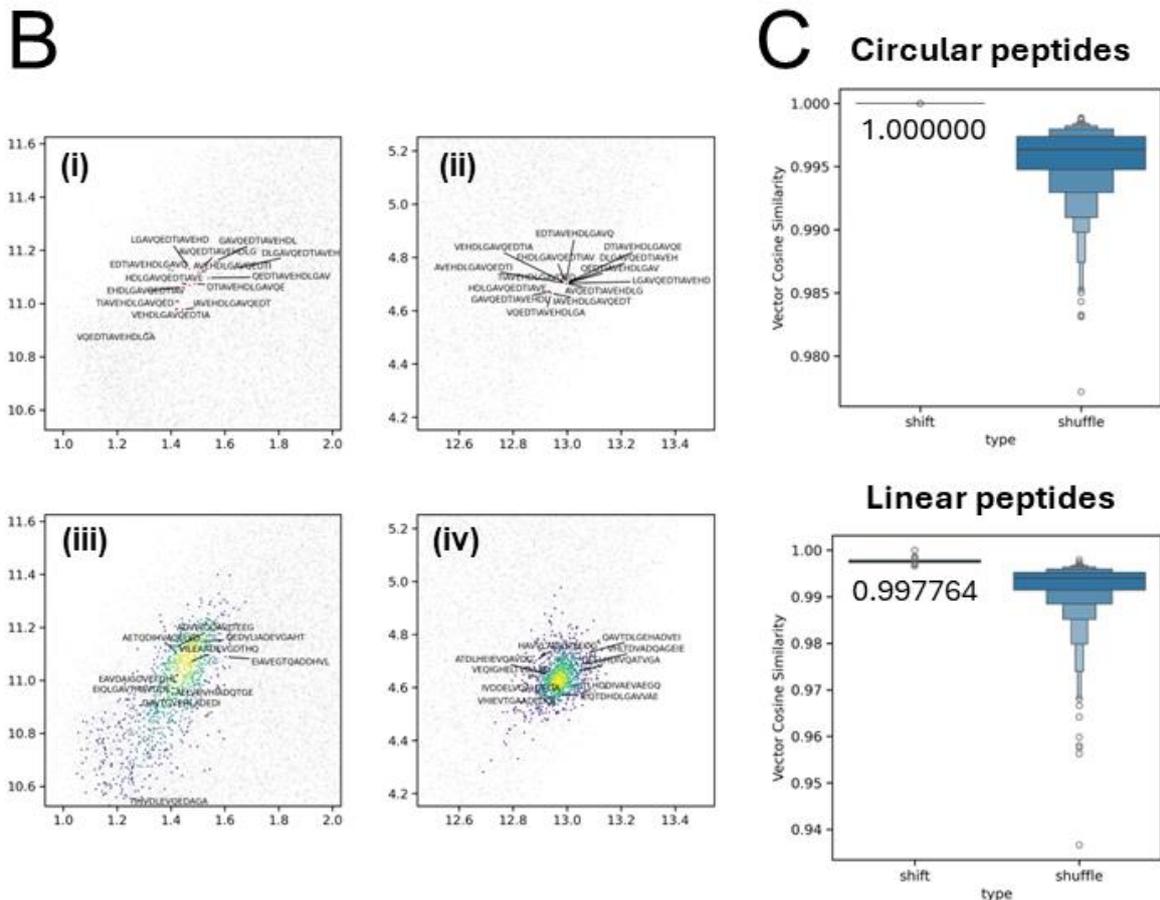
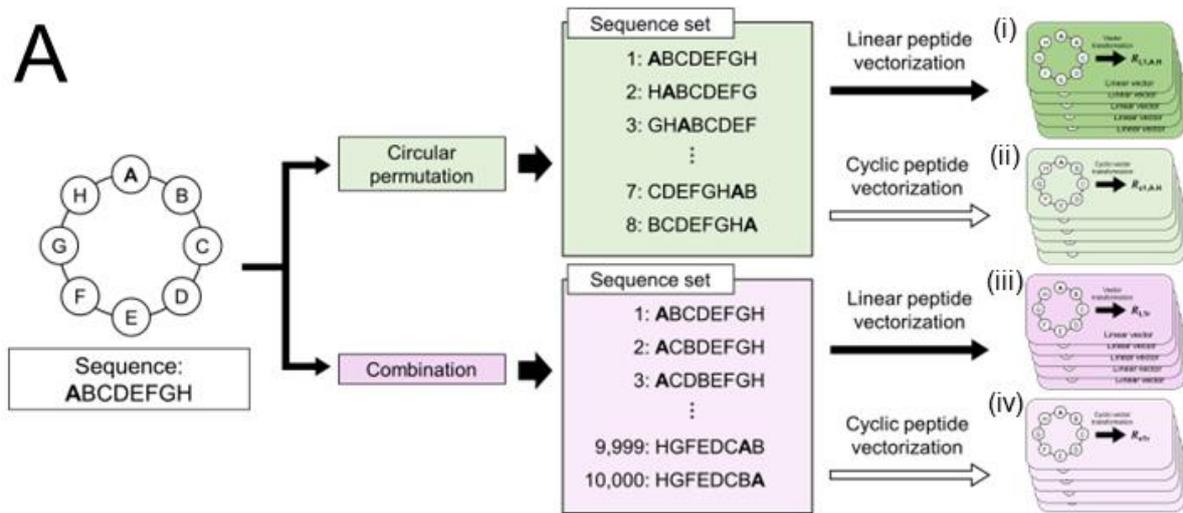
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12

## Supplementary information

### Supporting figures



**Supporting Figure 1 | UMAP projection of linear peptide embeddings.** Randomly generated 14-residue sequences were encoded using the ESM-2 protein language model. Unlike the cyclic vectorization strategy (Fig. 1A), the cyclic permutation averaging step was omitted (indicated by the crossed-out schematic), and sequences were processed as linear inputs. The resulting high-dimensional vectors were projected into a two-dimensional space using Uniform Manifold Approximation and Projection (UMAP). Density plots are displayed for datasets containing 1,000(upper), 10,000 (middle) and 300,000 (bottom) sequences.

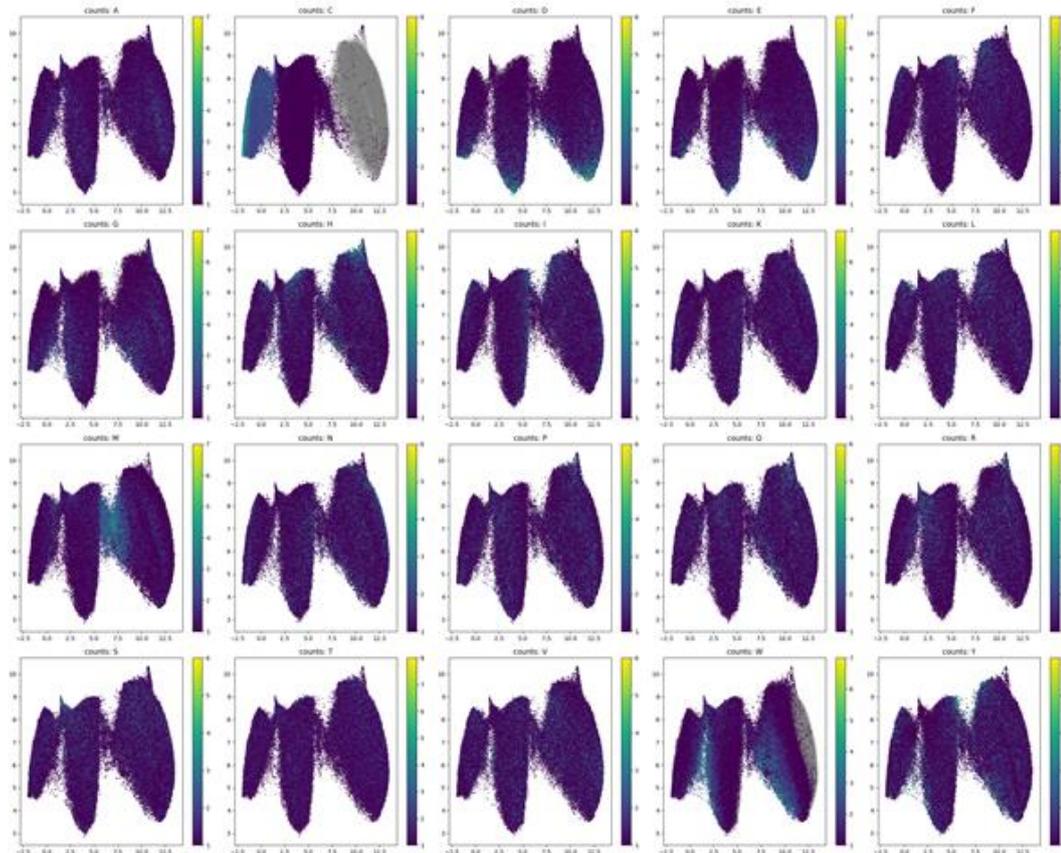


1  
2 **Supporting Figure 2 | Validation of permutation invariance and sequence specificity in cyclic peptide**  
3 **embedding.** (A) Experimental design for assessing vector robustness. A seed sequence (DTIAVEHDLGAVQE)  
4 was used to generate four distinct datasets to disentangle topological effects from compositional variance: linear  
5 and cyclic vectorization applied to (i, ii) 14 cyclic permutations (sequence shifts) and (iii, iv) 10,000 randomized  
6 permutations (shuffled combinations). (B) Visualization of vector convergence. UMAP projections of the four

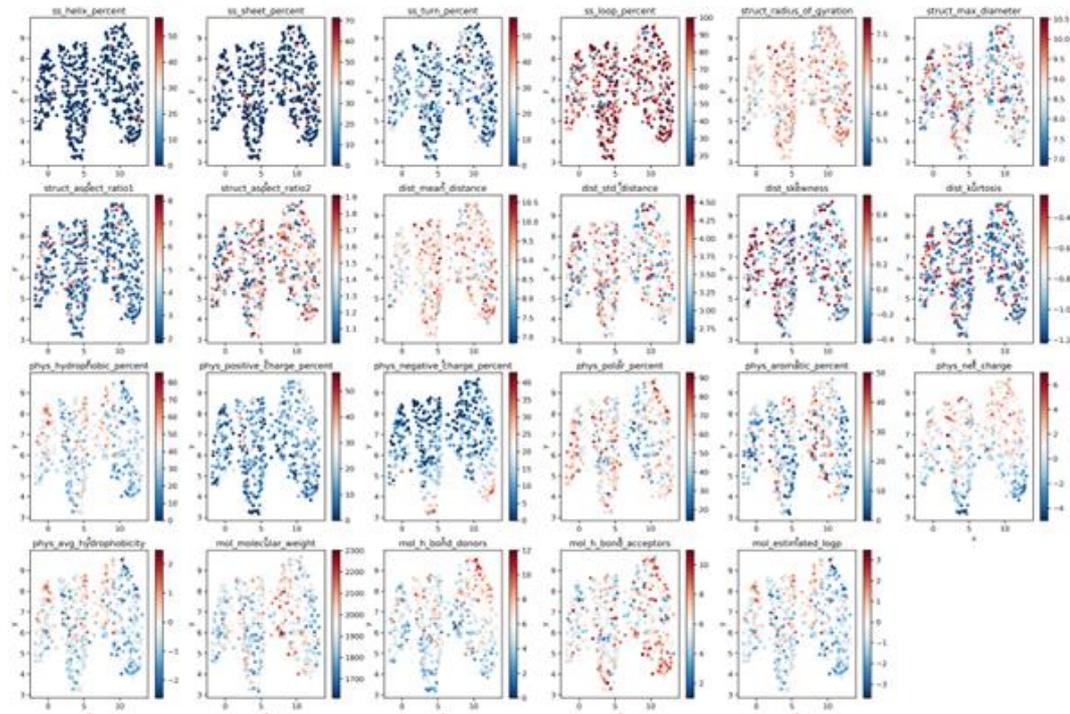
1 datasets. Notably, cyclic vectors derived from cyclic permutations (panel **ii**) converge to a single point, visually  
2 demonstrating that the embedding is independent of the starting residue. In contrast, linear vectors (panel **i**) exhibit  
3 slight spatial divergence due to N- and C-terminal edge effects. Both shuffled datasets (panels **iii** and **iv**) show  
4 broad spatial distribution, confirming that the embedding captures sequence order rather than mere amino acid  
5 composition. Since the peptide space is defined by the UMAP-based coordinate system, the x and y axis labels  
6 are omitted. **(C) Quantification of embedding stability.** Box plots of vector cosine similarities relative to the seed  
7 sequence. The cyclic vectorization method achieves perfect stability (similarity = 1.0000) across all cyclic  
8 permutations, effectively neutralizing the floating-point divergence observed in linear processing (similarity =  
9 0.9991). The marked variance in shuffled sequences further validates that the vector representation is strictly  
10 determined by the specific sequential arrangement of residues.

11

A



B



1

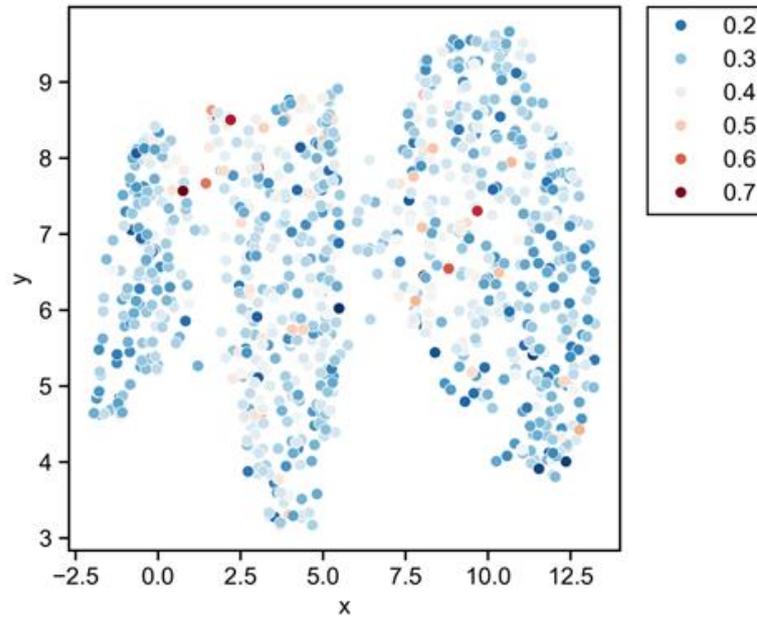
2 **Supporting Figure 3 | Deconvolution of compositional and structural determinants in the cyclic peptide**

3 **latent space. (A) Mapping of amino acid abundance distributions. The UMAP landscape is color-coded by the**

1 count of individual amino acid types within the constituent sequences. Cysteine content (panel 'counts: C') acts as  
2 the primary driver of the global topology, dictating the segmentation into three distinct clusters. Other residues,  
3 such as methionine (M) and tryptophan (W), also exhibit specific localized biases, highlighting that simple  
4 compositional variations significantly skew the spatial distribution. **(B)** Landscape of structural and  
5 physicochemical properties. To correlate the latent space with physical attributes, 3D structures for 3,000  
6 randomly sampled sequences were predicted using AfCycDesign [15]. Various metrics—including secondary  
7 structure content (e.g., helix/sheet percentage), geometric parameters (e.g., radius of gyration), and  
8 physicochemical nature (e.g., hydrophobicity, charge)—were projected onto the manifold. The resulting "mosaic-  
9 like" patterns demonstrate that these properties are not uniformly distributed but are instead organized into discrete  
10 regimes. This heterogeneity underscores that stochastic sequence generation fails to guarantee uniform sampling  
11 of the structural and physicochemical property space. Since the peptide space is defined by the UMAP-based  
12 coordinate system, the x and y axis labels are omitted.

13

14



1

2 **Supporting Figure 4 | Spatial localization of high-potential binder candidates. Visualization of design**

3 **performance landscapes.** The Loss values derived from EvoBind2 optimization targeting  $\beta 2m$  are mapped onto

4 the corresponding coordinates of the starting sequences in the peptide space. The heatmap color scale represents

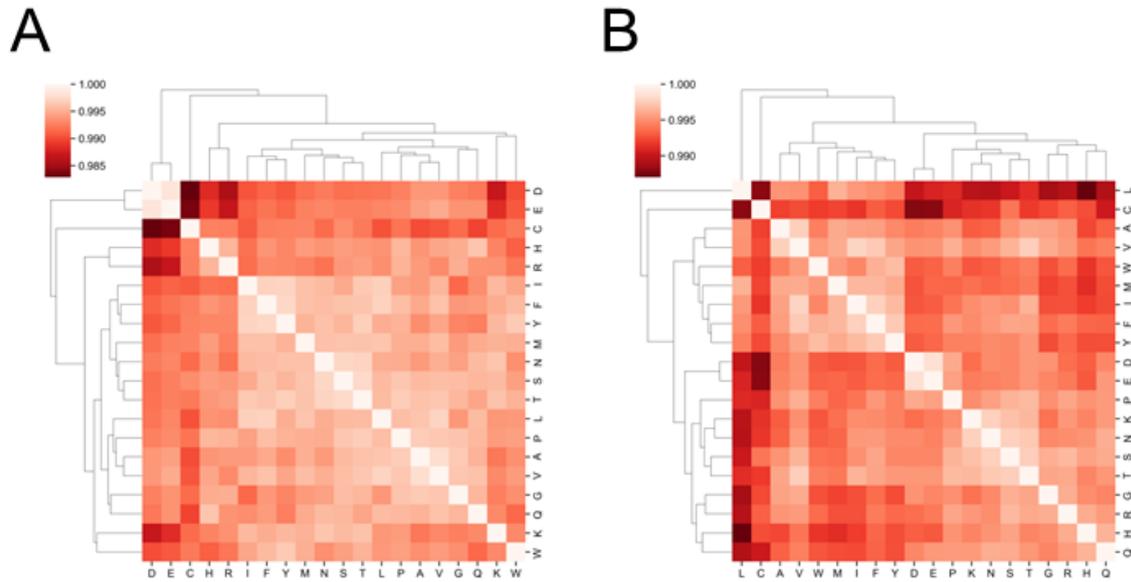
5 the raw Loss value, where specific points (e.g., darker hues) indicate sequences with superior predicted binding

6 properties (lower loss). Since the peptide space is defined by the UMAP-based coordinate system, the x and y axis

7 labels are omitted.

8

9



1

2 **Supporting Figure 5 | Hierarchical clustering of vector similarities for single-point mutations. (A, B)**

3 Pairwise cosine similarity matrices for peptide vectors derived from all 20 amino acid substitutions at position 2

4 (A) and position 6 (B) of the reference sequence. The heatmap color scale represents the degree of similarity

5 between mutant vectors, ranging from identity (white, 1.000) to lower similarity (dark red). Dendrograms along

6 the axes depict the hierarchical clustering of amino acid residues based on the calculated cosine

7 distances between their corresponding cyclic peptide vectors.

8

9