

Physics-based generative model of curvature sensing peptides; distinguishing sensors from binders

Niek van Hilten¹, Jeroen Methorst¹, Nino Verwei¹ and Herre Jelger Risselada^{1,2,3*}

¹Leiden Institute of Chemistry, Leiden University, Einsteinweg 55, Leiden, 2333 CC, The Netherlands.

²Department of Physics, Technical University Dortmund, Otto-Hahn-Strasse 4, Dortmund, 44227, Germany.

³Institute of Theoretical Physics, Georg-August-University Göttingen, Friedrich-Hund-Platz 1, Göttingen, 37077, Germany.

*Corresponding author(s). E-mail(s):

jelger.risselada@tu-dortmund.de;

Contributing authors: n.van.hilten@lic.leidenuniv.nl;

Abstract

Proteins can specifically bind to curved membranes through curvature-induced hydrophobic lipid packing defects. The chemical diversity among such curvature ‘sensors’ challenges our understanding of how they differ from general membrane ‘binders’, that bind without curvature selectivity. Here, we combine an evolutionary algorithm with coarse-grained molecular dynamics simulations (Evo-MD) to – for the first time – resolve the peptide sequences that optimally recognize the curvature of lipid membranes. We subsequently demonstrate how a synergy between Evo-MD and neural networks (NN) can enhance the identification and discovery of curvature sensing peptides and proteins. To this aim, we benchmark a physics-trained NN model against experimental data and show that we can correctly identify known ‘sensors’ and ‘binders’. We illustrate that sensing and binding are in fact phenomena that lie on the same thermodynamic continuum, with only subtle but explainable differences in membrane binding free energy, consistent with the serendipitous discovery of sensors.

Teaser

AI-based design helps explain curvature-selective membrane binding behavior.

1 Introduction

The recognition of curved regions of lipid bilayer membranes by proteins plays a key role in many biological processes, such as vesicular transport, fusion, and fission [1, 2]. This preferred binding to curved membranes is called curvature sensing and is driven by the outer leaflet of the curved bilayer membrane being stretched, which causes defects in the packing of the polar lipid head groups. Apolar amino acids of proteins can complement the now exposed hydrophobic tails within these lipid packing defects, negating their energetic penalty and resulting in a thermodynamic driving force (Fig. 1A).

Besides fundamental biological importance, curvature selectivity has been proposed as a potential avenue for the development of broad-spectrum antiviral peptides that leverage the difference in curvature between the membranes of small enveloped viruses and the essentially flat host cell membrane [3–6]. However, the extremely serendipitous discovery and resulting rarity of curvature selective peptides obstructs the utilization of state-of-the-art data-science driven generative models, like the recent work by IBM on the discovery of new antimicrobial peptides [7]. Consequently, an efficient computational strategy for accelerating the discovery of new curvature sensing peptides is still lacking.

Many natural curvature sensing proteins feature an amphipathic helix (AH). AHs have a polar face that interacts with the solvent and the lipid head groups and an apolar face that interacts with the hydrophobic lipid tails. Beyond this shared structural amphipathicity, the chemical composition of AHs is highly diverse. For example, the contrasting compositions of the amphipathic lipid packing sensing (ALPS) motif of the ArfGAP1 protein [8] and the AH of α -synuclein [9] (Fig. 1B-C) suggest that curvature sensing results from a delicate balance between the amino acid content on the apolar and polar sides of the helices [10, 11]. Moreover, and important to note, some AHs (like α -synuclein) have a positive net charge, providing additional selectivity for anionic liposomes specifically [12]. Taken together, the structural diversity among curvature sensors complicates reliable prediction of a given peptide’s sensing ability simply from sequence-based physicochemical descriptors, like mean hydrophobicity $\langle H \rangle$, hydrophobic moment μ_H [13], and net charge z .

Molecular dynamics (MD) simulations are a valuable asset in expanding our understanding of curvature sensing, since they can access the necessary molecular resolution that many experimental methods lack [14–16]. To reduce system size and, consequently, reduce the computational cost, curved membranes are often represented as stretched flat membranes in MD simulations [17, 18] (Fig. 1A), such that the lipid packing defects on the surface are similar and the consequent relative binding free energies correlate [19]. This approximation is valid since biologically relevant bilayer curvatures are negligible on the length scale of peptides (≈ 5 nm).

2 Results & Discussion

2.1 Designing the optimal curvature sensor

Evolutionary molecular dynamics (Evo-MD) is a physics-based inverse design method that embeds molecular dynamics (MD) simulations in a genetic algorithm (GA) framework [26]. GAs are inspired by Darwinian evolution and can serve as a powerful tool for optimization problems in large discrete search spaces, like the 20^L possible peptide sequences (for 20 natural amino acids and peptide length L). Starting from a, in our case random, initial subset ('population'), a GA iteratively (1) evaluates the desired property ('fitness') of the candidate solutions in the population, (2) selects the best candidates as the 'parents' for the next generation, and then (3) performs genetic operations, like cross-over recombination and random point mutations to (4) generate the next population (Fig. 2A). While evolution proceeds, the population's average fitness will increase until it converges to an optimum. To date, GAs have mainly been applied to peptide optimization problems that involve protein-peptide interactions and use fitness functions based on physicochemical descriptors or information from databases [20, 27–29]. In contrast, the fitness calculation in Evo-MD is based on ensemble averaging from (coarse-grained) MD simulation trajectories and is therefore completely data-independent. In this physics-based approach, experimental data contributes to solving the optimization problem via the parametrization of the force-field that is used in the simulations, the Martini model [30] in this study. Therefore, the main advantage is that Evo-MD will generate curvature sensing peptides without requiring any knowledge of existing curvature sensing peptides, of which too few examples exist to properly train a data-informed model. Additionally, as apposed to data-trained models, physics-based inverse design does not tend to generate molecules that are (too) similar to the input data [31]. In contrast, Evo-MD will search for a pre-defined thermodynamic optimum of sensing and generate 'optimal sequences' that actually may differ from the biological optimum due to additional evolutionary constraints imposed by nature's complexity (e.g. solubility, protein-protein interactions, trafficking).

The direction of simulated evolution by GAs is governed by the definition of the fitness function (the desired property). For the optimization of curvature sensing peptides, we aim to maximize the curvature sensing free energy $\Delta\Delta F$ that we can efficiently quantify using aforementioned mechanical free-energy method [19]. Our fitness function is the product of the $\Delta\Delta F$ value and a scaling factor c that equals 1 when located on the membrane surface and goes to 0 for transmembrane or soluble configurations (see SI). We emphasize that $\Delta\Delta F$ characterizes the relative affinity for lipid packing defects or equivalently positive leaflet curvature, analogous to the curvature-dependent binding constants (free energy of partitioning) measured in experimental model liposome assays [11, 23–25, 32–36] as well as measured differences in the concentration of peripheral membrane proteins due to curvature-driven sorting in micropipette aspiration assays [37, 38].

Since membrane-surface peptides are in most cases α -helical and the Martini force-field is unable to model protein folding events, we assume and fix helical secondary structure when generating the starting conformations for our peptides. To this end, and to reduce the search space, we excluded 10 amino acids with low α -helical propensities [39] (P, G, D, R, C, T, N, H, V, and I), whilst ensuring that every chemical subtype is represented in our final subset (comprising A, L, M, K, Q, E, W, S, Y, and F). We chose a fixed peptide length of 24 residues, which is in the typical range for curvature sensing peptides. Consequently, our search space contains 10^{24} peptide sequences.

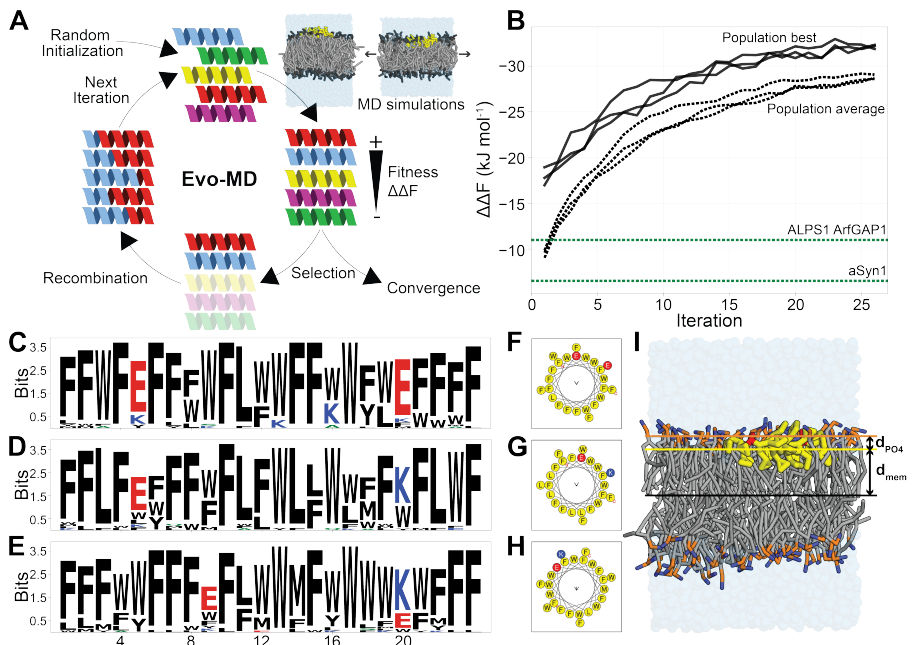


Fig. 2 **A)** Schematic representation of the Evo-MD process. Figure adapted from previous work [19]. **B)** Three independent replica Evo-MD runs show convergence within 25 iterations, as evident by the population best (solid lines) and population average (dashed lines). For comparison, the $\Delta\Delta F$ values for two known curvature sensing motifs (ALPS1 of ArfGAP1 [8] and α -synuclein [9]) are shown in green. **C-E)** Consensus sequence logos [40] for the best 36 sequences of the final population. **F-H)** The respective helical wheel representations [20] of the consensus sequences shown in C-E. **I)** Simulation snapshot of a consensus peptide (Fig. 2C, F) bound to a tensionless POPC membrane. Hydrophobic residues (F, W, L) are shown in yellow; E is shown in red. Phosphate (PO_4) and choline (NC_3) beads are shown in orange and blue, respectively. d_{mem} is the z -component of the center-of-mass distance between the membrane and the peptide. d_{PO_4} is the z -component of the center-of-mass distance between the PO_4 groups and the peptide.

In the three independent Evo-MD runs we performed (see SI for methodological details), we observed convergence within 25 iterations with the best candidates having a $\Delta\Delta F$ of around -32 kJ mol^{-1} (Fig. 2B). The consensus sequence logos [41] of the final generations show a strong enrichment of bulky

hydrophobic residues, mainly F and W (Fig. 2C-E and supplementary movies). We can understand this result by returning to our earlier statement that “a peptide that senses tension/curvature will also tend to induce tension/curvature”. Such induction of tension and concomitant membrane curvature in a membrane leaflet occurs via shallow insertions within the hydrophobic interior of the lipid membrane, i.e. the region directly below the head groups. Indeed, we observe that the optimal sequences – the point of maximum leaflet tension generation with respect to the helix’ central axis – is characterized by an insertion of 1.69 ± 0.05 nm from the bilayer center in a tensionless membrane (d_{mem} in Fig. 2I), or alternatively 0.24 ± 0.05 nm below the average position of the phosphate groups (d_{PO_4} in Fig. 2I). This is in quantitative agreement with predictions from membrane elastic theory, which suggest an optimal insertion of about 1.7 nm from the membrane center plane [42]. Furthermore, the bulkier the peptide is, and thus the larger its excluded volume and effective helical radius, the more pronounced the induced leaflet tension will be.

Besides the abundant hydrophobic residues, we observed that the solutions of all three Evo-MD runs feature two charged residues (E or K, see Fig. 2C-E). This numerical conservation of two charged amino acids suggests that this is the bare minimum of polar content that is necessary to maintain a surface orientation for such hydrophobic peptides (i.e. scaling factor $c \rightarrow 1$). The sign of the charge appears to be irrelevant for the zwitterionic POPC membranes we used here. Also, the exact position of these residues seems arbitrary, as long as the two charged residues end up on the same side of the helix in the folded conformation (Fig. 2F-H).

The fact that all three randomly initiated Evo-MD runs produced peptide sequences with identical physical characteristics within the same number of iterations strongly suggests that this is indeed the global and not a local optimum. To probe the effect of only using the 10 most helix-prone amino acids, we performed an additional Evo-MD run with all 20 natural residues included. This, again, yielded peptides with the same physical characteristics, but showing slower convergence (40 iterations) and higher diversity due to the vastly increased search space (see SI).

What this simulated evolution shows is that the GA has successfully selected a key aspect in curvature sensing, insertion of hydrophobic residues [23, 36], which is then maximally amplified and exploited until the fitness converges. To such extent even, that the optimal ‘sensor’ is so hydrophobic that it would likely stick to any membrane, regardless of curvature, thus being classified as a ‘binder’ instead. What is immediately clear is that our optimized peptides strongly differ from the naturally evolved optima (e.g. the ALPS motif and α -synuclein), both in terms of $\Delta\Delta F$ (Fig. 2B) and in their chemical compositions (Fig. 1B-C and Fig. 2F-H). Thus, our physics-based inverse design indicates that the distinction between curvature sensors and membrane binders can be considered as a continuum with a soft, subtle threshold at a relative binding free energy that is much lower than the theoretical optimum.

Biologically, the differences between the simulated optimum and naturally evolved peptides can be explained by considering the many boundary conditions imposed by the complex environment of *in vivo* systems. One of the most obvious and fundamental requirements is that proteins should be soluble in physiological buffer. The extremely hydrophobic GA-generated optima clearly fail this criterion and will readily aggregate and precipitate out of the solution. Also, since curvature sensing implies tension generation, peptides with a high $\Delta\Delta F$ could harm the integrity of the membranes they adhere to. Therefore, an evolutionary pressure to limit this potency must exist.

2.2 A neural network model to predict curvature sensing

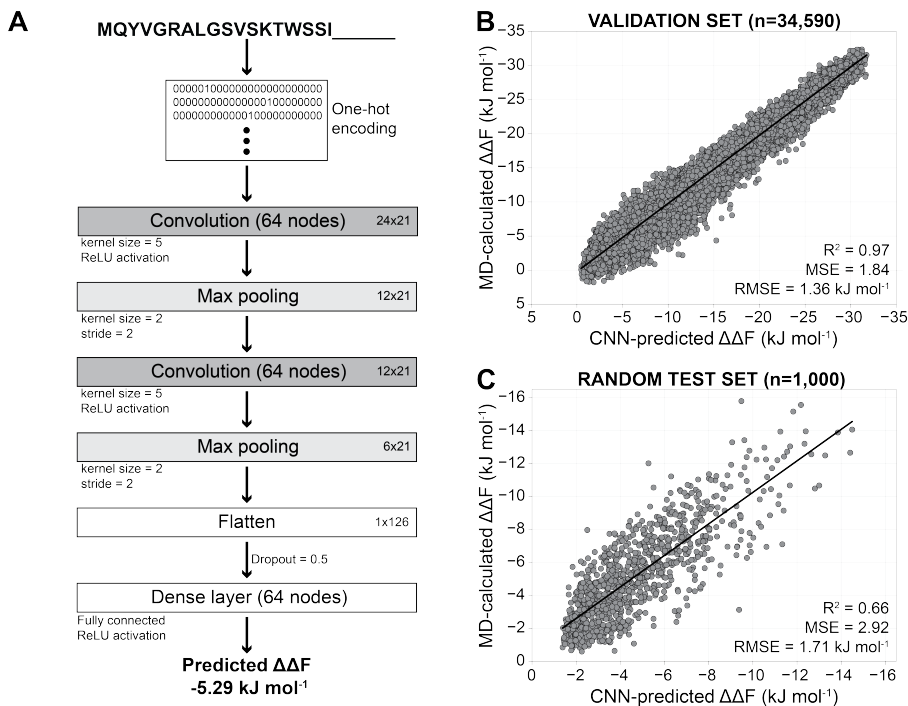
As a valuable byproduct of the iterative optimization process by Evo-MD, we obtained a large database of $\approx 54,000$ unique sequences (all 24 residues long) and their respective sensing free energies ($\Delta\Delta F$) as calculated by MD simulations. With this wealth of data, we set out to train a convolutional neural network (CNN) that is able to predict curvature sensing ability from peptide sequence information only.

To enable the model to handle peptides shorter than 24 amino acids as well, we split the sequences in the original data set at a random position, such that the resulting two fragments were at least 7 residues long. Next, since $\Delta\Delta F$ depends linearly on length (see SI), we interpolated the $\Delta\Delta F$ values for the split sequences, hereby tripling the data set to $\approx 138,000$ sequences (after discarding duplicate fragments). We refrained from extrapolating to sequences longer than 24 residues, since this would require additional assumptions on amino acid composition and – potentially – involve more complex tertiary structures that are inaccurately modeled by the Martini force-field. A detailed description of the final training data is included in the SI.

As described previously in the context of activity prediction of helical antimicrobial peptides [43], we used one-hot encoded and zero-padded representations for the input sequences. These are then fed to two consecutive convolutional layers with max pooling, followed by a fully connected layer and a single output neuron to translate the connection weights into a float value: the predicted $\Delta\Delta F$ (Fig. 3A, see SI for details on the optimization of the architecture and hyperparameters).

During the CNN training, minimization of the mean squared error (MSE) converged after 18 epochs (see SI) to a MSE of 1.84 for the validation set (25% random sample from the full data). For this validation set, we achieved excellent correlation ($R^2 = 0.97$) between the predicted and MD-calculated $\Delta\Delta F$'s (Fig. 3B). However, because sequences from the late iterations of the same Evo-MD run can be highly similar, the validation and training sets are arguably not fully independent. Therefore, to ultimately test our model, we predicted the $\Delta\Delta F$ for 1,000 randomly generated sequences (between 7 and 24 residues long) that were not part of the training data and obtained a MSE of 2.92 and an R^2 -value of 0.66 when comparing the predicted values to the $\Delta\Delta F$ calculated by MD simulations (Fig. 3C).

8 Generative model of curvature sensing peptides



The trained neural network and all data sets are accessible at github.com/nvanhiltten/CNN_curvature_sensing. Please note that the model should only be used for sequences between 7 and 24 amino acids long, and that it assumes α -helical folding (as we did in the training data). Based on the performance of our model on the randomly generated test set, the root-mean-square error (RMSE) of its predictions is $\sqrt{MSE} = 1.71$ kJ mol⁻¹, which is comparable to the typical errors obtained when calculating $\Delta\Delta F$ by MD simulation (e.g. compare the error bars of Fig. 4A and 4B).

2.3 Distinguishing sensing from binding

Now, with the MD quantification and CNN tools in hand, we can return to the key question posed in the introduction, namely: which, if any, characteristics can help us distinguish curvature sensors from binders, and what can relative binding free energies, like $\Delta\Delta F$, teach us in this regard? To address this question, we composed a benchmark set of natural curvature sensing peptides (see SI), also including mutated variants that were empirically categorized as ‘non-binders’ (i.e. no affinity for any membrane) or ‘binders’ (i.e. binding to membranes without curvature specificity). We should acknowledge the expert help of Prof. Bruno Antonny and Dr. Romain Gautier in composing this list.

To fairly compare the sequences, we propose two correction factors to obtain an adjusted relative binding free energy $\Delta\Delta F_{adj}$. First, we linearly extrapolate the $\Delta\Delta F$ values of shorter peptides to their corresponding free energies if they were 24 residues long ($\Delta\Delta F_{L=24}$, see SI). Second, we realized that many of the peptides are cationic to improve interaction with (curved) anionic membranes that are abundant in nature. Since our MD simulations were performed with neutral POPC membranes, we hypothesized that the relative binding free energies would in these cases be underestimated and thus require a correction term $c_z z$ to account for this:

$$\Delta\Delta F_{adj} = \Delta\Delta F_{L=24} + c_z z \quad (1)$$

To determine the magnitude of c_z (the relative free energy contribution per unit charge z), we performed additional MD simulations with anionic membranes (75% POPC, 25% POPG) and indeed found elevated relative binding free energies, especially for the cationic peptides (see SI). From the average difference between the $\Delta\Delta F$'s calculated on the different membranes, we obtained $c_z = -0.93 \pm 0.89$ kJ mol⁻¹ per unit charge.

We calculated this $\Delta\Delta F_{adj}$ for our peptides of interest with the CNN model (Fig. 4A) and MD simulations (Fig. 4B). When ranking the peptides accordingly, we find that we can roughly reproduce the empirical qualification of ‘non-binders’ at the lower end and ‘binders’ at the higher end of the ranking, and ‘sensors’ in the middle. Notably, the differences between the values for subsequently ranked peptides are rather small, often below 1 kJ mol⁻¹. These findings are in line with our hypothesis that the thermodynamics of binding and sensing are subtle energetic transitions on a continuous scale.

Interestingly, we find that the CNN-predicted ranking is in better agreement with the experimental trends than the results from MD simulations (Fig. 4A-B). We speculate that this is likely due to a smoothening effect, i.e. predictions of the MD simulations for individual peptide sequences are fully independent whereas the CNN introduces effective correlations between sequences. Because the CNN is trained on the ensemble averaged values from many thousands of independent MD trajectories for different peptide sequences, we argue that disturbances (‘noise’) in chemical space (point mutations) and in the molecular dynamics itself (limited sampling) are therefore smoothened out to such extent that the experimental trends are more robustly reproduced. Hence, we use the CNN-predicted $\Delta\Delta F_{adj}$ values for the remainder of this paper.

2.4 Thermodynamic model of the sensing→binding transition

Along the lines of our current definitions, every peptide undergoing hydrophobically driven membrane binding is able to sense positive membrane curvature due to the increase in surface hydrophobicity upon bending. In other words, positive curvature enhances hydrophobically driven membrane binding. Empirically, however, a peptide is only classified as a curvature sensor if it *only*

significantly binds to membranes characterized by a high positive curvature. In this section, we will argue that the empirical classification of sensors versus binders can be intuitively understood from the population statistics of a two-state partition function.

Herein, we define the following two states; (1) state m : the peptide is bound to the membrane, (2) state s : the peptide is in solution. We define the partitioning free energy difference between the two states as ΔF_{sm} : the free energy of membrane binding minus the free energy of solvation with respect to the peptide in the gas phase. We calculated ΔF_{sm} by thermodynamic integration (see SI) for a non-binder, a sensor, a binder, and the extremely hydrophobic Evo-MD optimum (Fig. 2D, F), and found that it linearly relates to $\Delta\Delta F_{adj}$ (Fig. 4C). The reason for this strong linear correlation is that both membrane binding (ΔF_{sm}) and curvature sensing ($\Delta\Delta F_{adj}$) are driven by hydrophobic interactions.

At thermal equilibrium, the relative probability for a peptide to bind to a membrane P_m is Boltzmann distributed. Consequently, P_m is given by:

$$P_m = \frac{1}{1 + \frac{V_s}{A_m} e^{\Delta F_{sm}}} \quad (2)$$

Eq. 2 resembles the so-called Fermi-Dirac function, in which V_s represents the normalized volume of accessible solvent (proportional to the number of realizations in solution) and A_m represents the normalized membrane area (proportional to the number of membrane binding realizations). If $\frac{V_s}{A_m} = 1$, this function features a sharp but continuous transition at $\Delta F_{sm} = 0$ (dashed curve in Fig. 4E), i.e. the point where membrane binding and peptide solubility are precisely in balance ($P_m = 0.5$). However, the number of realizations in solution is expected to be much larger than the number of realizations associated with membrane binding, $\frac{V_s}{A_m} \gg 1$, at typical lipid concentrations. Consequently, the transition point shifts to the right, i.e. peptide-membrane binding is significantly favoured over peptide solvation at the transition point $P_m = 0.5$. The steep nature of this ‘switch function’ strikingly explains *why* empirically classified sensing behavior switches to binding behavior based on only subtle differences in (relative) membrane binding free energy, like we observed in this work (Fig. 4A).

To finalize our model, we estimate the prefactor $\frac{V_s}{A_m}$. The accessible membrane area in one liter of solution is $A = \frac{1}{2}cN_{av}A_{lip}$, with c being the lipid concentration, N_{av} the Avogadro constant, and A_{lip} the area per lipid. The characteristic surface area of a helical peptide is roughly $A_p = 5 \times 1 \text{ nm}^2$ and equivalently its volume in solution is $V_p = 5 \times 1 \times 1 \text{ nm}^3$. Taking $c = 2 \text{ mM}$ (a typical value in the middle of the concentration range used in experiments [11, 23–25, 32–36]) and $A_{lip} = 0.64 \text{ nm}^2$ we obtain $\frac{V_s}{A_m} = \frac{V/V_p}{A/A_p} = 2.6 \cdot 10^3$.

When we plug this number into eq. 2, we find that the free energy of membrane binding outperforms the free energy of solvation by $-19.5 \text{ kJ mol}^{-1}$ at the transition point $P_m = 0.5$ (Fig. 4E), i.e. the transition point is associated with

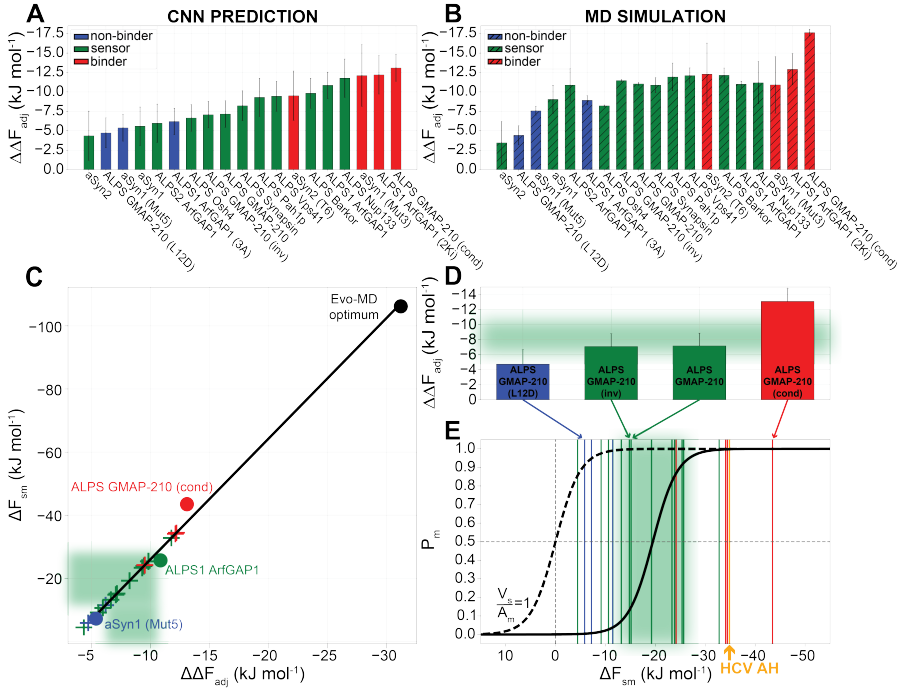


Fig. 4 **A-B)** Length- and charge-adjusted relative binding free energy ($\Delta\Delta F_{adj}$, see eq. 1) for 19 curvature sensors or derivatives (see SI) predicted with the CNN model (A) and calculated from MD simulations (B). Non-binders are shown in blue. Sensors are shown in green. Binders are shown in red. **C)** Linear correlation between CNN-predicted $\Delta\Delta F_{adj}$ and the membrane binding free energy (ΔF_{sm}). Circles indicate peptides for which ΔF_{sm} was calculated by thermodynamic integration (see SI). The Evo-MD optimum (black) is the sequence in Fig. 2D. For the remaining peptides (crosses), ΔF_{sm} was derived from the linear fit $\Delta F_{sm} = 3.83 \cdot \Delta\Delta F_{adj} + 12.27$. **D)** Highlighted CNN-predicted $\Delta\Delta F_{adj}$ values for ALPS GMAP-210 variants. **E)** The probability of a peptide to bind a membrane P_m as a function of the membrane binding free energy ΔF_{sm} (eq. 2). The dashed curve is for $\frac{V_s}{A_m} = 1$, which is shifted to the solid curve for $\frac{V_s}{A_m} = 2.6 \cdot 10^3$, as derived in the main text. In C-E, the green area indicates the likely regime ($0.05 \leq P_m \leq 0.95$) where peptides are empirically classified as sensors. The orange line indicates the curvature specific antiviral peptide HCV AH.

favourable but relatively weak membrane binding. Because of the sharp transition behavior, ‘sensors’ are envisioned as peptides with a ΔF_{sm} value near the transition point ($0.05 \leq P_m \leq 0.95$, the green area in Fig. 4E). According to our model, a peptide with $P_m < 0.05$ would be classified as a ‘non-binder’, and a peptide with $P_m > 0.95$ would be deemed a ‘binder’. When overlaying the values for the 19 benchmark peptides we introduced earlier (vertical lines in Fig. 4E, see SI), we find that 3/3 non-binders, 7/11 sensors, and 3/4 binders are classified correctly, i.e. in agreement with the original experiments. Such categorization would have been impossible with crude physicochemical descriptors like mean hydrophobicity $\langle H \rangle$ or hydrophobic moment μ_H despite them (weakly) correlating with ΔF_{sm} (see SI).

We defined the sensing regime rather generously ($0.05 \leq P_m \leq 0.95$) to accommodate for the fact that both the empirical classification and the positioning of the Fermi-Dirac curve are sensitive to the lipid concentration c , which can differ significantly between experiments. Moreover, imperfect helical folding and differing lipid compositions, cell systems, and read-out types can complicate the direct comparison between empirical classifications and our computational models. In terms of membrane binding free energy ΔF_{sm} , the lower and upper boundaries of sensing ($0.05 \leq P_m \leq 0.95$) correspond to -12.2 and -26.8 kJ mol $^{-1}$, respectively (green area in Fig. 4C, E). In other words, when the work required to pull a peptide from a tensionless membrane is between those values, it can be classified as a sensor. In terms of the length- and charge-adjusted relative binding free energy ($\Delta\Delta F_{adj}$), i.e. the preference for highly curved membranes (≈ 25 nm vesicles [19]), sensors fall between -6.4 and -10.2 kJ mol $^{-1}$ (green area in Fig. 4C-D).

As an example, we highlighted four variants of ALPS GMAP-210 [24, 25] for which the relative differences in experimental sensing/binding behavior were correctly captured (Fig. 4D). This is a striking example because the sequences are so similar: the only difference between the original sensor ALPS GMAP-210 and the non-binding variant ALPS GMAP-210 (L12D) is a single point mutation (L \rightarrow D) that disturbs the peptide's hydrophobic face. As can be expected, and in line with the experimental findings, the inverse sequence ALPS GMAP-210 (inv) scores the same as the original peptide and is thus categorized as a sensor as well. Finally, the condensed version ALPS GMAP-210 (cond) was correctly identified as a binder. With this highlighted example, we demonstrate that our method and thermodynamic model can pick out features as subtle as single point mutations, resolve the resulting differences in relative free energy and correctly categorize the consequent sensing behavior.

Finally, we also included the antiviral peptide HCV AH, that specifically ruptures vesicles with a high curvature (e.g. small enveloped viruses) [3, 4] (orange line in Fig. 4E). To date, HCV AH is the only example of a clinically relevant curvature selective antiviral peptide [5]. When we plug in the previously calculated free energy value [19] for HCV AH, we find that this peptide falls into the 'binder' regime. This is consistent with evidence that the vesicle size specificity of this peptide is due to curvature specific pore formation and not to curvature specific binding (i.e. curvature sensing) [44]. After all, subtle binding may not be optimal for pore formation, i.e. the subsequent induction of tension should be sufficient to rupture the membrane.

Hence, we argue that the most promising range to find potent curvature specific antiviral agents is therefore near the transition zone between sensing and binding (i.e., $P_m \rightarrow 1$), since these peptides (1) may still benefit from some 'curvature sensing' (predominant binding to higher curvatures), (2) pack a larger punch than biological sensors in terms of meeting the tension induction threshold necessary to deform/perforate viral membranes, but are (3) not so potent that they also rupture the host cell membrane. The latter is helped by the fact that the host membrane is likely more resilient to the disruptive

actions of peptides than the viral membrane due to membrane stabilizing proteins and active feedback mechanisms. This means that a considerable part of the selectivity of membrane-targeting drugs is likely due to a difference in drug resistance (membrane resilience) rather than actual differential binding.

3 Conclusions

We have illustrated a striking example of the utility of a physics-based generative model (Evo-MD) to explore and simultaneously rationalize the mechanisms of *how* peptides sense membrane curvature. Initially, we set out to optimize curvature sensing by resolving the sequence that maximizes the *relative* affinity for lipid packing defects. Instead, we ended up with the optimal ‘binder’. This finding led to the important realization that curvature sensing and membrane binding are phenomena that lie on the same thermodynamic continuum (Fig. 4E).

Naturally evolved curvature sensors, such as the ALPS motif and α -synuclein, are chemically diverse but turn out to be remarkably similar in terms of partitioning free energies, which explains their functional similarities. In this work, we – for the first time – described the thermodynamic regime that defines the curvature sensing behavior of peptides. Given how narrow this energetic ‘sensing window’ is, it is unsurprising that the discovery and design of curvature sensors has, to date, been rather serendipitous. Having identified this ‘window of opportunity’ in terms of relative binding free energy can facilitate the discovery of novel curvature sensing peptides since we now know where to look for them.

The existence of the here-resolved thermodynamic sensing regime can also be intuitively understood from an evolutionary biological perspective. Hence, curvature sensing motifs within naturally evolved proteins must fulfil the following two criteria: they should (1) predominantly bind to curved membranes and (2) conserve the structural integrity of the membranes they adhere to. The here-observed linear correlation between sensing and overall membrane binding (see Fig. 4C) dictates that these criteria are only met in the weak binding regime. These arguments are all in full agreement with the earlier hypothesis that curvature sensing in nature is a subtle balance between overall membrane binding and specific curvature recognition [2]. Also, in this weak binding regime, the small leaflet strain induced by peptides is able to facilitate a largely inert and thus biologically functional sensing phenotype that does not easily lead to membrane rupture/deformation.

Importantly, we demonstrated a fruitful synergy between a physics-based generative model (Evo-MD) and a convolutional neural network, which not only dramatically accelerated the high-throughput evaluation of peptide sequences, but also improved the accuracy of prediction compared to the original molecular simulations. It is important to stress the key role of Evo-MD,

in that it yields sequences over the whole range of $\Delta\Delta F$ by gradually maximizing the relevant chemical property in a well spaced manner, in our example even up to the thermodynamic optimum. Using such data to train a neural network model has the important advantage that it encompasses the full thermodynamic range of possibilities over a vast search space of 20^{24} sequences, whereas a data set of natural peptides (if available in the first place) would be strongly constrained to a certain biologically feasible regime that only comprises peptides with highly similar physicochemical characteristics. We argue that training data generated with Evo-MD can therefore substantially improve both the applicability domain as well as the accuracy of neural network models, despite many of the generated sequences not being necessarily biologically relevant. This principle is equivalent to the fitting of an unknown function to data points that are well spaced over the whole range of the applicability domain versus data points that are only clustered within a narrow window. Particularly, precise knowledge of the maxima (and minima) of a function – which a physics-based optimization resolves – will benefit the quality of a fit or model, also within the biologically relevant domain of the search space. We postulate that a subsequent restriction of the search space within or near the here-resolved sensing regime can enhance the discovery of novel curvature sensing motifs in natural proteins, as well as their *de novo* generation.

Finally, we envision an important potential application in the computational design of peptide sensors that recognize membranes with other aberrant characteristics, such as a distinct lipid composition (e.g. bacteria and cancer cells). Since (selective) membrane binding results in the generation of leaflet tension, membrane binding peptides have an inherent membrane destabilizing propensity as well as the ability to lower the energetic cost of the highly curved interface of toroidal pores. This is particularly the case for the hydrophobically driven membrane binding peptides we discussed in this work. The simultaneous encoding of selective membrane binding plus an active drug mode, such as the induction of membrane lysis, is therefore a realistic avenue to explore further. An important advantage of physics-based generative models over existing data-science based generative models herein is their unique ability to systematically explore the drug therapeutic potential of distinct relative binding ($\Delta\Delta F$) regimes by restricting the generation of peptide sequences within pre-defined boundary values of $\Delta\Delta F$, for example, via the straightforward introduction of a bias/constraint to the fitness function. This can enable the targeted exploration of different ‘windows of opportunity’ similar to aiming a gun at different targets.

Supplementary information. Methodological details and additional analyses are included in the supplementary information (SI). We also provide the evolution of consensus sequences as supplementary movies.

Acknowledgments. We thank Prof. Bruno Antony and Dr. Romain Gauthier for insightful discussions and their expert help in selecting relevant curvature sensors, non-binders, and binders to test our theory and methods.

The Dutch Research Organization NWO (Snellius@Surfsara) and the HLRN Göttingen/Berlin are acknowledged for the provided computational resources. This work was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy - EXC 2033 - 390677874 - RESOLV. We thank the NWO Vidi scheme (project number 723.016.005), and the DFG (grant number RI2791/2-1) for funding.

Author contributions. N.v.H. and H.J.R. designed the research. J.M. wrote the Evo-MD code. N.v.H implemented Evo-MD modules specific to the curvature sensing problem. N.v.H. and N.V. developed, trained, and optimized the neural network model. N.v.H. and H.J.R. wrote the manuscript.

Data availability. The Evo-MD code and all related files are available at github.com/nvanhilten/Evo-MD_curvature_sensing.

The data sets, trained CNN model, and related scripts are available at github.com/nvanhilten/CNN_curvature_sensing.

References

- [1] N. S. Hatzakis *et al.*, *Nat. Chem. Biol.* **5**(11), 835 (2009).
- [2] B. Antonny, *Annu. Rev. Biochem.* **80**, 101 (2011).
- [3] N.-J. Cho *et al.*, *ACS Chem. Biol.* **4**(12), 1061 (2009).
- [4] J. A. Jackman, G. H. Zan, V. P. Zhdanov, N.-J. Cho, *J. Phys. Chem. B.* **117**(50), 16117 (2013).
- [5] J. A. Jackman *et al.*, *Nat. Mater.* **17**(11), 971 (2018).
- [6] B. K. Yoon, W.-Y. Jeon, T. N. Sut, N.-J. Cho, J. A. Jackman, *ACS Nano* **15**(1), 125 (2021).
- [7] P. Das *et al.*, *Nat. Biomed. Eng.* **5**(6), 613 (2021).
- [8] J. Bigay, J. F. Casella, G. Drin, B. Mesmin, B. Antonny, *EMBO J.* **24**(13), 2244 (2005).
- [9] M. B. Jensen *et al.*, *J. Biol. Chem.* **286**(49), 42603 (2011).
- [10] G. Drin and B. Antonny, *FEBS Lett.* **584**(9), 1840 (2010).
- [11] I. M. Pranke *et al.*, *J. Cell Biol.* **194**(1), 89 (2011).
- [12] W. S. Davidson, A. Jonas, D. F. Clayton, J. M. George, *J. Biol. Chem.* **273**(16), 9443 (1998).
- [13] D. Eisenberg, R. M. Weiss, T. C. Terwilliger, *Nature* **299**(5881), 371 (1982).
- [14] P. González-Rubio, R. Gautier, C. Etchebest, P. F. J. Fuchs, *Biochim. Biophys. Acta Biomembr.* **1808**(9), 2119 (2011).
- [15] S. Vanni, H. Hirose, H. Barelli, B. Antonny, R. Gautier, *Nature Comm.* **5**, 4916 (2014).
- [16] K. S. Stroh and H. J. Risselada, *J. Chem. Theory Comput.* **17**(8), 5276 (2021).
- [17] N. van Hilten, K. S. Stroh, H. J. Risselada, *Front. Physiol.* **11**, 250 (2020).
- [18] C.-P. Chng, N.-J. Cho, K. J. Hsia, C. Huang, *Langmuir* **37**(45), 13390 (2021).
- [19] N. van Hilten, K. S. Stroh, H. J. Risselada, *J. Chem. Theory Comput.* **18**(7), 4503 (2022).

- [20] R. Gautier, D. Douguet, B. Antonny, G. Drin, *Bioinformatics* **24**(18), 2101 (2008).
- [21] M. G. Ford *et al.*, *Nature* **419**(6905), 361 (2002).
- [22] C. Löw *et al.*, *Biophys. J.* **95**(9), 4315 (2008).
- [23] G. Drin *et al.*, *Nat. Struct. Mol. Biol.* **14**(2), 138 (2007).
- [24] G. Drin, V. Morello, J.-F. Casella, P. Gounon, B. Antonny, *Science* **320**(5876), 670 (2008).
- [25] M. Magdeleine *et al.*, *eLife* **5**, e16988 (2016).
- [26] J. Methorst, N. van Hilten, H. J. Risselada, *bioRxiv* (2021).
- [27] I. Belda, X. Llorà, E. Giralt, *Soft Comput.* **10**(4), 295 (2006).
- [28] B. Knapp, V. Giczi, R. Ribarics, W. Schreiner, *BMC Bioinformatics* **12**(1), 241 (2011).
- [29] S. J. Barigye, J. M. García de la Vega, Y. Perez-Castillo, J. A. Castillo-Garit, *Future Med. Chem.* **13**(11), 993 (2021).
- [30] P. C. T. Souza *et al.*, *Nature Methods* **18**(4), 382 (2021).
- [31] B. Sanchez-Lengeling and A. Aspuru-Guzik, *Science* **361**(6400), 360 (2018).
- [32] B. Mesmin *et al.*, *Biochemistry* **46**(7), 1779 (2007).
- [33] W. Fan, A. Nassiri, Q. Zhong, *Proc. Natl. Acad. Sci.* **108**(19), 7769 (2011).
- [34] M. Cabrera *et al.*, *J. Cell Biol.* **191**(4), 845 (2010).
- [35] E. Karanasios, G.-S. Han, Z. Xu, G. M. Carman, S. Siniossoglou, *Proc. Natl. Acad. Sci.* **107**(41), 17539 (2010).
- [36] S. Vanni *et al.*, *Biophys. J.* **104**(3), 575 (2013).
- [37] A. Tian and T. Baumgart, *Biophys. J.* **96**(7), 2676 (2009).
- [38] B. Sorre *et al.*, *Proc. Natl. Acad. Sci.* **106**(14), 5622 (2009).
- [39] C. N. Pace and J. M. Scholtz, *Biophys. J.* **75**(1), 422 (1998).
- [40] A. Tareen and J. B. Kinney, *Bioinformatics* **36**(7), 2272 (2019).

- [41] T. D. Schneider and R. M. Stephens, *Nucleic Acids Res.* **18**(20), 6097 (1990).
- [42] F. Campelo, H. T. McMahon, M. M. Kozlov, *Biophys. J.* **95**(5), 2325 (2008).
- [43] J. Witten and Z. Witten, *bioRxiv* (2019).
- [44] S. R. Tabaei, M. Rabe, V. P. Zhdanov, N.-J. Cho, F. Höök, *Nano Lett.* **12**(11), 5719 (2012).