

## RESEARCH ARTICLE

# Differentiable molecular simulation can learn all the parameters in a coarse-grained force field for proteins

Joe G. Greener<sup>1</sup>\*, David T. Jones

Department of Computer Science, University College London, London, United Kingdom

\* [j.greener@ucl.ac.uk](mailto:j.greener@ucl.ac.uk)

## Abstract

Finding optimal parameters for force fields used in molecular simulation is a challenging and time-consuming task, partly due to the difficulty of tuning multiple parameters at once. Automatic differentiation presents a general solution: run a simulation, obtain gradients of a loss function with respect to all the parameters, and use these to improve the force field. This approach takes advantage of the deep learning revolution whilst retaining the interpretability and efficiency of existing force fields. We demonstrate that this is possible by parameterising a simple coarse-grained force field for proteins, based on training simulations of up to 2,000 steps learning to keep the native structure stable. The learned potential matches chemical knowledge and PDB data, can fold and reproduce the dynamics of small proteins, and shows ability in protein design and model scoring applications. Problems in applying differentiable molecular simulation to all-atom models of proteins are discussed along with possible solutions and the variety of available loss functions. The learned potential, simulation scripts and training code are made available at <https://github.com/psipred/cgdms>.

## OPEN ACCESS

**Citation:** Greener JG, Jones DT (2021) Differentiable molecular simulation can learn all the parameters in a coarse-grained force field for proteins. *PLoS ONE* 16(9): e0256990. <https://doi.org/10.1371/journal.pone.0256990>

**Editor:** Yang Zhang, University of Michigan, UNITED STATES

**Received:** April 19, 2021

**Accepted:** August 19, 2021

**Published:** September 2, 2021

**Peer Review History:** PLOS recognizes the benefits of transparency in the peer review process; therefore, we enable the publication of all of the content of peer review and author responses alongside final, published articles. The editorial history of this article is available here: <https://doi.org/10.1371/journal.pone.0256990>

**Copyright:** © 2021 Greener, Jones. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** The learned potential, simulation scripts and training code are made available under a permissive license at <https://github.com/psipred/cgdms>.

## Introduction

Molecular simulation has been successful in making predictions and understanding experimental data [1, 2]. Treating the system with an appropriate level of complexity, usually all-atom molecular mechanics or residue-level coarse-graining when simulating proteins, is necessary to access the timescales required for the property under investigation [3, 4]. It is generally agreed that force fields are not optimally parameterised [5, 6], for example significant effort has gone into modifying a handful of parameters in standard force fields to better represent both ordered and disordered proteins [7–9]. Such efforts have improved the force fields without changing their functional form, an attractive proposition when the alternative is adding complexity that restricts the timescales available for study.

Meanwhile, deep learning has had a major impact on many areas of biology, achieving state of the art performance in fields such as protein structure prediction [10]. A number of groups have applied these advances to molecular simulations [11–13] including learning coarse-grained potentials [14–18], learning quantum mechanical potentials [19–22], improving

**Funding:** This work was supported by the European Research Council (<https://erc.europa.eu>) Advanced Grant “ProCovar” (project ID 695558) awarded to DTJ. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing interests:** The authors have declared that no competing interests exist.

sampling [23, 24], and improving atom typing [25]. Whilst promising, many of these approaches show limited success when used on systems they were not trained on. Methods trained on trajectory data also suffer from a lack of standardised simulations across many systems due to the high cost of obtaining such trajectories. Other groups have used inventive machine learning strategies for end-to-end protein structure prediction [26] and to score static conformations [27, 28].

The idea of differentiating through the numerical solution of Newton’s equations of motion to obtain gradients that can be used to improve a learned force field is called differentiable molecular simulation (DMS). It has been discussed [29], but is yet to produce generally useful force fields. Conceptually the idea is related to recent work on neural differential equations [30–32]. The idea is appealing because a large number of parameters can be improved at once, rather than the small numbers currently modified. The gradients are also exact, at least with respect to the numerical integration and the loss function used. Only recently has the hardware and software been available to run such simulations. The variety of available loss functions and ability to calculate exact gradients for all parameters suggest DMS could be the next step in the steady improvement of force fields [33].

Previous studies utilising DMS have trained neural network potentials, ranging from graph neural networks [29] to ambitious protein folding simulators running Langevin dynamics [34]. Whilst training neural networks may be an effective solution, as discussed above there is much room for improvement in existing force fields. It makes sense to try and improve these where possible rather than moving to a new functional form, which has the additional advantage of retaining the physical interpretability of conventional force fields. Neural networks are also slower to run than existing force fields, meaning that to avoid reducing available simulation time the network has to be trained to jump multiple time steps, a rather different problem to learning the instantaneous potential. This work is also influenced by a number of studies that compare native and training ensembles to improve force fields for protein folding using maximum likelihood, contrastive divergence and related approaches [35–45]. These methods are able to modify many parameters at once, but they generally involve comparing conformations rather than obtaining gradients of some loss function through the simulation and are hence limited in the properties they can target.

In this study we use automatic differentiation (AD) [46], the procedure used to train neural networks where it is called backpropagation, to learn all the parameters from scratch in a simple coarse-grained force field for proteins. This learned potential matches potentials derived from chemical knowledge and Protein Data Bank (PDB) statistics, reproduces native flexibility when used in simulation, is able to fold a set of small proteins not used for training, and shows promise for protein design and model scoring applications. It adds to existing coarse-grained and statistical potentials used for simulation [47–50] and model scoring approaches [51–53]. More broadly, it points to DMS as a useful technology falling under the banner of differentiable programming [54], an expansion of the principles of deep learning to the concept of taking gradients through arbitrary algorithms and utilising the known structure of the system under study [55, 56].

## Results

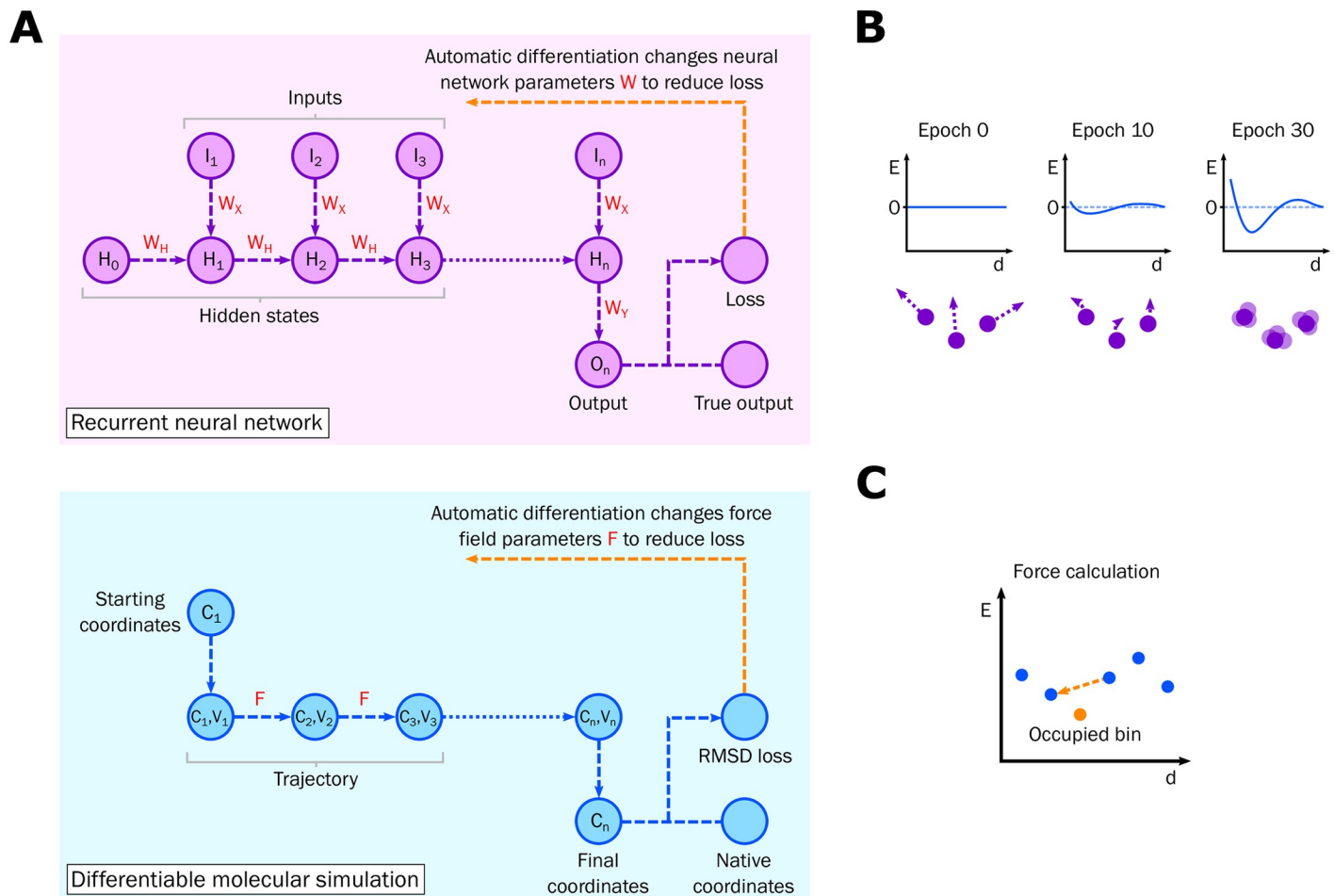
### Differentiable molecular simulation

In this study a coarse-grained potential for proteins is learned in which a protein is represented by 4 point particles per residue (N, C $\alpha$ , C and sidechain centroid) with no explicit solvent. The potential consists of 3 components: pairwise distance potentials (including covalent bonds), bond angle potentials, and torsion angle potentials associated with the predicted secondary

structure type of a residue. Overall there are 29,360 individual potentials and 4.1m learnable parameters. The high number of parameters indicates some redundancy due to the nature of the potentials, but also demonstrates that DMS can be used to learn a large number of parameters at once. See the methods for further details on the functional form of the potentials, how forces are calculated and how the model is trained. Proteins presented in the results do not have homologs present in the training set and single sequence secondary structure prediction is used throughout the study, so the results are not due to overtraining or direct learning of evolutionary information. Assessing the method on proteins not used for training distinguishes this approach from many approaches used to date for machine learning of molecular simulations.

During training, proteins are simulated using the velocity Verlet integrator in the NVE ensemble, i.e. with no thermostat. The starting conformation is the native structure and up to 2,000 steps are run. Due to parameters shared across each time step this can be thought of as analogous to a recurrent neural network (RNN) running on a sequence of length 2,000, as shown in Fig 1A. In particular, implementing the simulation in the AD framework PyTorch [57] allows the gradients of a given loss function with respect to each learned parameter to be calculated, allowing an optimiser to change the parameters to reduce the loss function. However, the learned parameters are those of the force field rather than the weights and biases of a standard neural network.  $\log(1 + R_f)$  is used as the loss function, where  $R_f$  is the root-mean-square deviation (RMSD) across all atoms in the coarse-grained model between the conformation at the end of the simulation and the native structure. At the start of learning the potentials are flat, the forces are zero and the proteins distort according to the randomised starting velocities, as shown in Fig 1B. Over the course of training, the  $R_f$  values decrease as the potential learns to stabilise the proteins in the training set. Previous studies have adopted a similar approach of minimising the final RMSD but used random sampling rather than exact gradients [58]. The training and validation  $R_f$  values throughout training are shown in S1 Fig in S1 File.

After training on a dataset of 2,004 diverse proteins up to 100 residues long the potentials resemble those derived from chemical knowledge and PDB statistics. Due to the coarse-grained nature of the simulation the energy values, along with other properties such as the time step and the temperature used later in the thermostat, cannot be assigned standard units. As shown in Fig 2A, covalent bond distance potentials have strong minima at the correct distance and steep barriers preventing steric clashing. The steep drops at the edges are an artifact of training and do not affect simulations or energy scoring, since these values are never occupied when using the trained potential. Bond angle potentials have minima at the true values with a few degrees of tolerance allowed either side, as shown in Fig 2B. Same residue C $\alpha$ -side-chain distance potentials indicate that different rotamer conformations have been learned; Fig 2C shows that the energy minima for isoleucine agree with the two minima found in the PDB distance distributions, which correspond to different rotamers. Other cases showing different rotamer conformations include glutamic acid, glutamine, lysine, methionine and tryptophan. These can be seen in the complete set of such potentials shown in S2 Fig in S1 File. The torsion angle potentials show different preferences for residues predicted as  $\alpha$ -helical,  $\beta$ -sheet and coiled, see Fig 2D, and these agree with the true Ramachandran distributions. The glycine torsion angle potentials indicate lower energy regions in the positive  $\phi$  space and the proline potentials display a minimum at 0° for the  $\omega$  angle corresponding to cis-proline (data not shown). The potentials most important for the tertiary structure are the general distance potentials. As shown in Fig 2E these match potentials of mean force (PMFs) derived from the PDB in many cases, with minima for many pairs around 6 Å driving hydrophobic packing. The steep energy barriers to steric clashing do not generally extend below 4 Å because



**Fig 1. Differentiable molecular simulation.** (A) The analogy between a RNN and DMS. Learnable parameters are shown in red. The same parameters are used at each step. There are many variants of RNNs; the architecture shown here has a single output for a variable length input, which could for example represent sentiment classification of a series of input words. (B) Learning the potential. A representation of a component of the potential is shown with energy  $E$  plotted against inter-atomic distance  $d$ . At the start of training (epoch 0) the potential is flat and the atoms deform according to their starting velocities. During training the potential learns to stabilise the native structures of the training set. (C) Force calculation from the potential. Adjacent bins to the occupied bin are used to derive the force using finite differences. In this case the force acts to reduce the distance  $d$ . See the methods for more details.

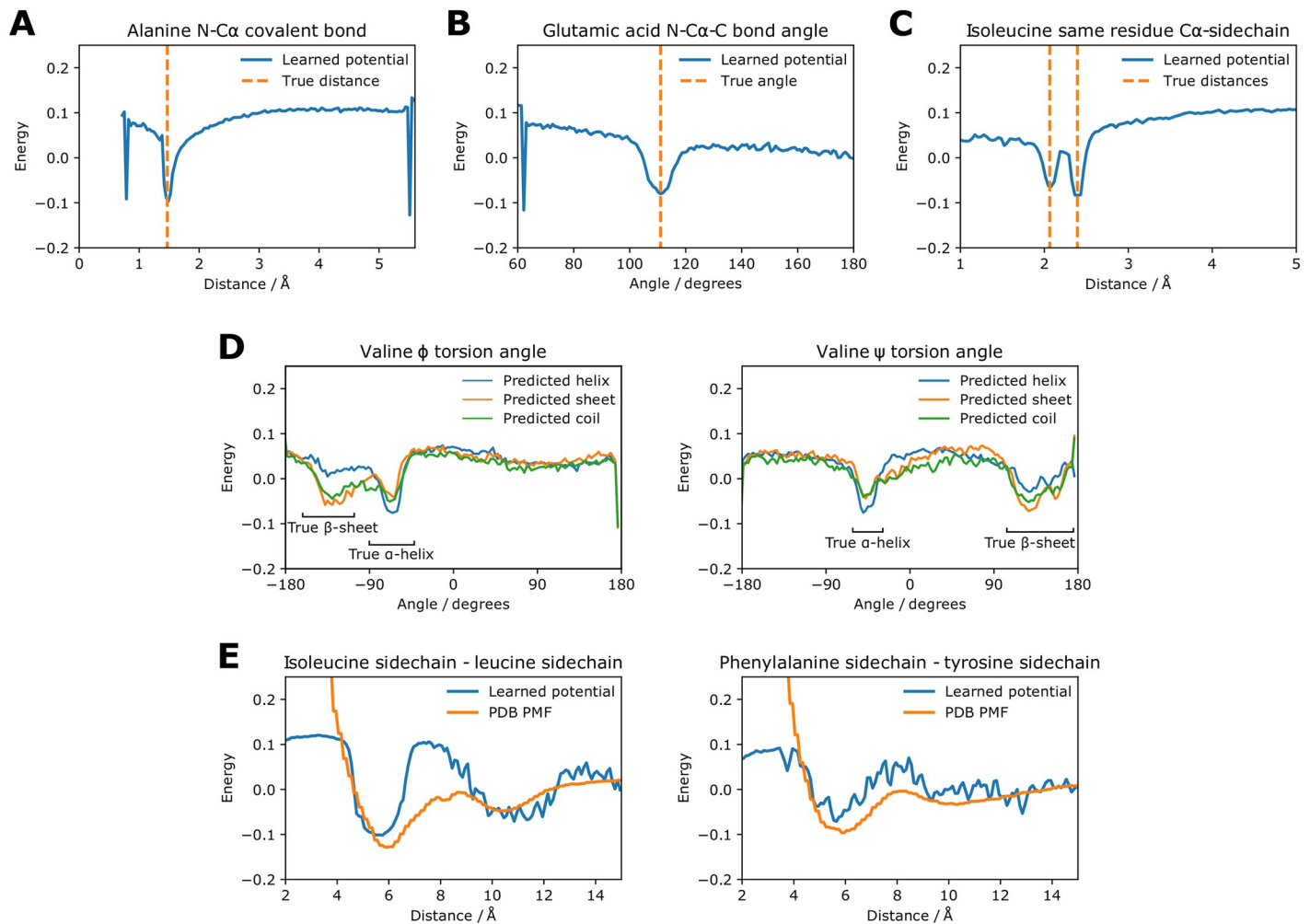
<https://doi.org/10.1371/journal.pone.0256990.g001>

extending these is not required to improve  $R_f$  during training. Steric clashing is not seen during simulations provided that a suitable time step is used. The complete set of potentials for general sidechain-sidechain distances are shown in S3 Fig in S1 File.

Despite being trained only to minimise the final RMSD across the protein, we find that the learned potential shows local detail. The potential energy when modifying the  $\phi$  and  $\psi$  torsion angles of alanine dipeptide is shown in Fig 3A. Shown in Fig 3B is the free energy calculated from an all-atom simulation in Wang et al. 2019 [15]. The learned potential matches the major low energy conformations of the all-atom model.

## Protein structure and dynamics

A learned potential can be used to run a simulation of arbitrary length since gradients are not recorded. Here we study four small, fast-folding proteins investigated with molecular simulation previously [8, 59, 60]. They all have NMR ensembles or crystal structures available from experiments [61, 62]. Details on all proteins presented in the results are in S1 Table in S1 File.

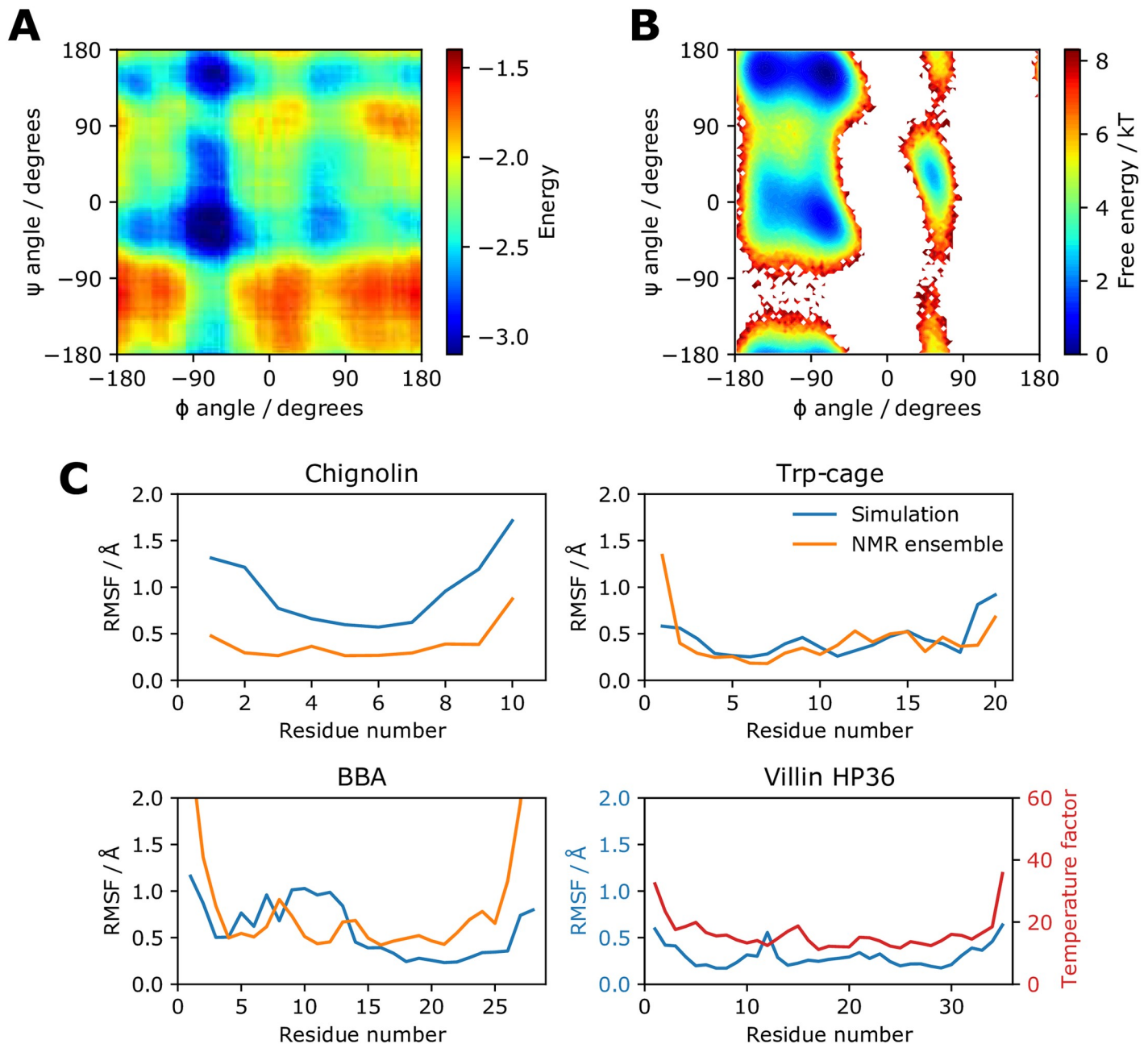


**Fig 2. Components of the learned coarse-grained potential compared to true values.** The energy scale is the same for each plot. Each individual potential consists of discrete energy values for 140 evenly-spaced bins. (A) Distance potential between N and C $\alpha$  in the same alanine residue, i.e. a covalent bond. (B) Bond angle potential for the N-C $\alpha$ -C angle in glutamic acid. (C) Distance potential between C $\alpha$  and the sidechain centroid on the same residue in isoleucine. (D) Torsion angle potentials in valine. There are different potentials for residues predicted as  $\alpha$ -helical,  $\beta$ -sheet and coiled. The true ranges of  $\alpha$ -helices and  $\beta$ -sheets are shown. (E) Distance potentials between sidechains for two pairs of amino acids. A PMF calculated from the PDB is also shown for comparison.

<https://doi.org/10.1371/journal.pone.0256990.g002>

First we investigate whether the learned potential can keep the native structures stable and reproduce residue-level flexibility. As shown in Fig 4B, the native structures are stable under simulation in the NVT ensemble using the Andersen thermostat, with C $\alpha$  RMSDs less than 4  $\text{\AA}$  in all cases. Fig 3C shows that the root-mean-square fluctuation (RMSF) of the C $\alpha$  atom of each residue over the simulation generally matches that of the native NMR ensembles, with an expected increase in flexibility for terminal residues. Chignolin displays more flexibility under simulation than in the NMR ensemble, likely due to the lack of explicit hydrogen bonding in the coarse-grained model to keep the  $\beta$  turn structure rigid. For villin HP36, we see general agreement between the residue RMSF and the crystal temperature factors.

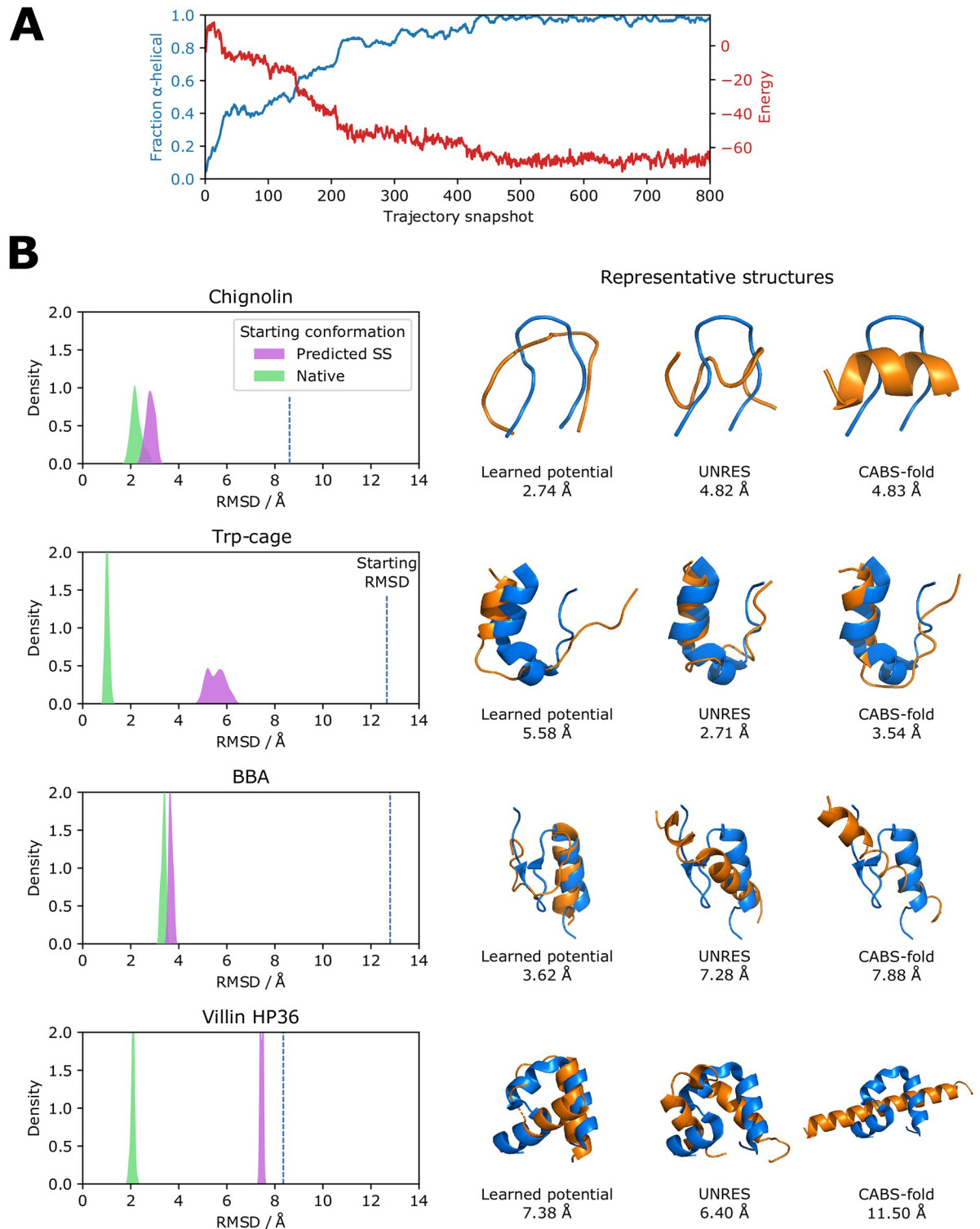
Next, we ask whether we can fold small proteins and peptides in the NVT ensemble. We find that the (AAQAA) $_3$  repeat peptide folds into an  $\alpha$ -helix over 12m steps when started from a random conformation, matching its experimental behaviour at physiological temperatures. This is accompanied by a reduction in the energy under the learned potential and is shown in



**Fig 3. Protein energy and dynamics.** (A) The energy of alanine dipeptide in the learned potential for different  $\phi$  and  $\psi$  torsion angles. Conformations are generated with PeptideBuilder at intervals of  $2^\circ$  and scored. (B) Free energy for different  $\phi$  and  $\psi$  torsion angles for all-atom alanine dipeptide. The energy is obtained from direct histogram estimation from all-atom simulations. The data is taken from Wang et al. 2019 [15] and provided by the authors. (C) C $\alpha$  atom RMSF values for simulations of 6m steps using the learned potential starting from the native structure. An initial burn-in period of 6m steps was discarded. The C $\alpha$  atom RMSF values from the NMR ensembles are also shown, or the C $\alpha$  atom crystal temperature factors in the case of villin HP36.

<https://doi.org/10.1371/journal.pone.0256990.g003>

**Fig 4A.** We do find that longer sequences are able to form the correct secondary structure, and often the correct tertiary structure given enough simulation time. However to better explore tertiary structure formation with available compute resources we started the proteins from extended conformations containing predicted secondary structure, i.e.  $\alpha$ -helical  $\phi/\psi$  angles for predicted  $\alpha$ -helical residues and extended  $\phi/\psi$  angles for residues predicted  $\beta$ -sheet or coiled.



**Fig 4. Folding small proteins and peptides.** (A) The  $(AAQAA)_3$  repeat peptide folds from a random starting conformation into an  $\alpha$ -helix over 12m steps. The  $\alpha$ -helical fraction and energy in the learned potential are shown with a snapshot taken every 15,000 steps. (B) C $\alpha$  RMSD distributions from simulations of 3m steps for 4 proteins starting from predicted secondary structure and native conformations. An initial burn-in period of 9m steps was discarded. The starting C $\alpha$  RMSD for the predicted secondary structure conformation is also shown. A representative structure found with MDAnalysis is shown (orange) along with models generated from the web servers of UNRES and CABS-fold. The C $\alpha$  RMSD to the native structure (blue) is given for each model.

<https://doi.org/10.1371/journal.pone.0256990.g004>

The C $\alpha$  RMSD distributions of the trajectories after a burn-in period, the starting C $\alpha$  RMSDs and representative models from the trajectories are shown in Fig 4B. All the proteins fold to a native-like structure over 12m steps, taking about 36 hours on a single graphics processing unit (GPU) or about 3 times as long on the central processing unit (CPU). Chignolin adopts approximately the correct  $\beta$  turn structure; the minimum sampled C $\alpha$  RMSD of 2.15 Å is comparable to the C $\alpha$  RMSD of 1.82 Å between the first model in the NMR ensemble and the crystal structure. Trp-cage lacks native helix formation in the middle of the protein but the overall shape is correct. BBA forms a native-like structure with a minimum C $\alpha$  RMSD during the simulation of 3.47 Å. For villin HP36 the N-terminal helix faces the wrong direction but the location of the turns and the rest of the structure is correct. This indicates a low energy basin in which the topological mirror structure is found, a problem that can be overcome with a higher temperature (see below) or enhanced sampling.

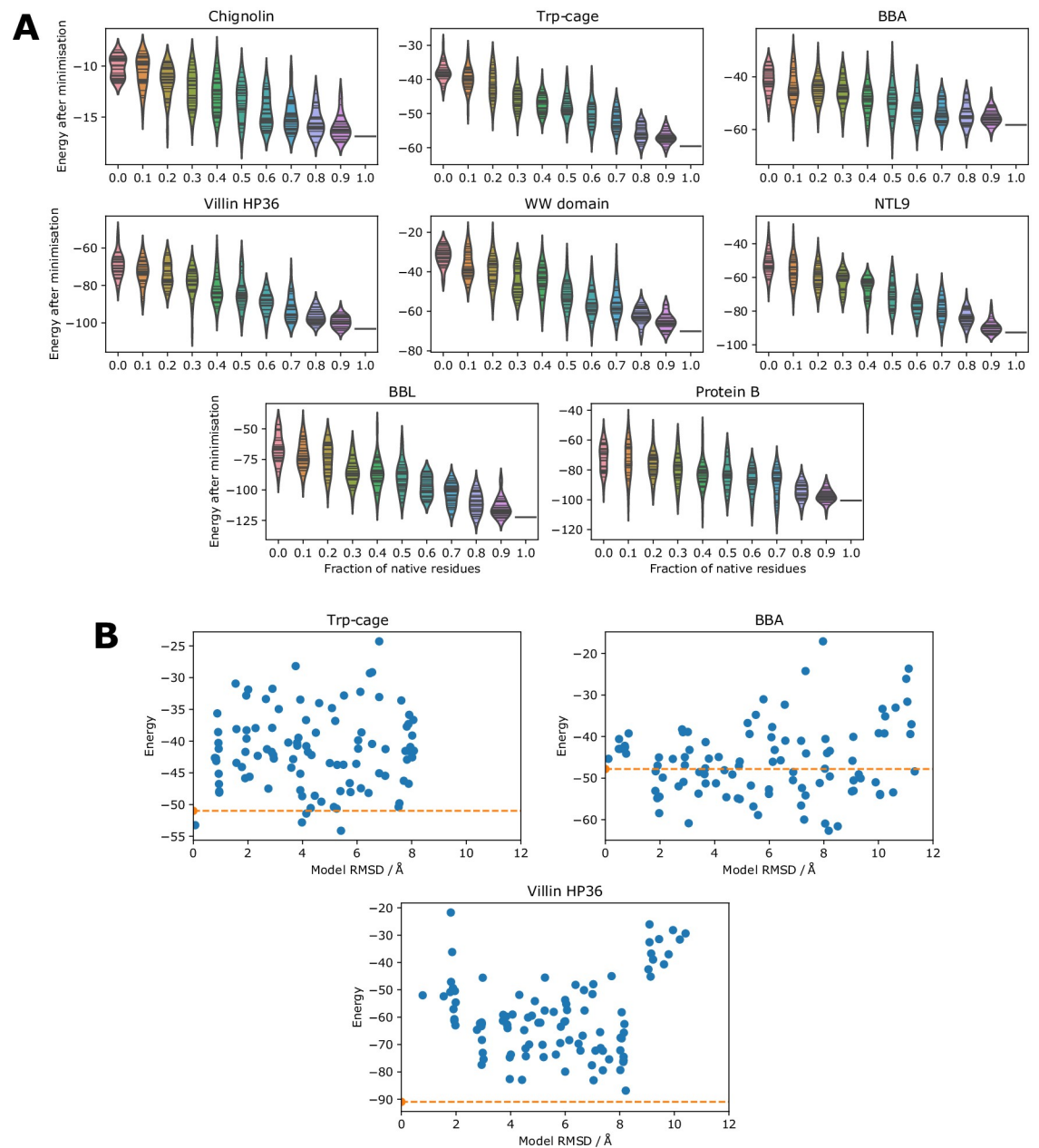
Models are also shown from web servers implementing two popular methods that carry out coarse-grained protein folding: UNRES [47, 63], which has two interacting sites per protein, and CABS-fold [48, 64], which uses a lattice model. Both make use of predicted secondary structure, and are given the same single sequence prediction used here. Performance of the learned potential is comparable to these established methods across the proteins tested, with the learned potential able to break secondary structure elements and add turns in the correct locations. The comparison to these methods is not exact; they provide less compute resources on their servers than the 36 hours of GPU time per protein used here, whereas this method does not employ the replica exchange algorithms used to enhance sampling in UNRES and CABS-fold. Enhanced sampling approaches to predict the structures of larger proteins with the learned potential is a topic of future work. We note that larger proteins remain close to their native structure with low energy when simulated, suggesting that the native structure is in an energy minimum that can be accessed with appropriate sampling. Another point of note is that the same simulation parameters are used when simulating all four proteins (temperature 0.015, coupling constant 25). More accurate models can likely be obtained by optimising these for each protein but we did not want to risk overfitting. The (AAQAA)<sub>3</sub> repeat peptide helix formation simulations were, however, carried out at a higher temperature (temperature 0.022, coupling constant 100) to faster explore the conformational space needed to form the  $\alpha$ -helix. We do notice some success in folding from an extended chain with these higher temperature parameters. For example, villin HP36 reaches a minimum C $\alpha$  RMSD of 4.19 Å over 12m steps, with the orientation of all the helices correct.

## Protein design and model scoring

In order to see whether native sequences are optimal for native structures in the learned potential, we thread sequences with varying fractions of native amino acids onto the native backbone and calculate the energy. Since long simulations are not required for these tests we also used four more proteins investigated with similar methods previously [59, 60]. Non-native residues are drawn randomly from the background distribution of amino acids in the PDB, i.e. leucine is more likely to be chosen (9.6% chance) than tryptophan (1.2% chance). As shown in Fig 5A we find that an increased fraction of native amino acids gives a lower energy, with the native sequence lower in energy than most 90% native sequences and considerably lower in energy than less native sequences.

Next we used the potential for fixed backbone design to see if designed sequences match the native sequence. We start from a random sequence and make mutations. At each trial a residue is mutated and the mutation is accepted or rejected based on the energy change when threaded onto the native structure and the distance through the trial process (see the methods).





**Fig 5. Protein design and model scoring.** (A) The energy of sequences with varying fractions of native residues threaded onto the native structure. 30 sequences are scored for each fraction. A short minimisation of 100 steps is carried out for each sequence and the energy is recorded at the end of the minimisation. (B) The energy of well-packed decoys generated using 3DRobot is compared to the C $\alpha$  RMSD of the decoys (blue dots) and the energy of the native structure (orange line). Chignolin was too small to run 3DRobot on.

<https://doi.org/10.1371/journal.pone.0256990.g005>

After 2,000 trials some of the designed sequences match the native sequences, particularly in the core regions with inter-residue interactions. The native fraction for each protein is shown in Table 1 along with analogous results from EvoEF2 [65], which uses an energy function consisting of nine terms and is developed specifically for protein design applications. The results from the potential presented here are comparable to those from EvoEF2, with our potential performing better on the smaller proteins and worse on the larger proteins. Training the

**Table 1. Fixed backbone design.**

Protein	EvoEF2 fraction of native residues	This work fraction of native residues
Chignolin	0.300	0.400
Trp-cage	0.300	0.550
BBA	0.143	0.179
Villin HP36	0.229	0.143
WW domain	0.303	0.152
NTL9	0.308	0.231
BBL	0.234	0.149
Protein B	0.106	0.085

The fraction of residues that match the native residue at the end of the design process is given for this work and for EvoEF2.

<https://doi.org/10.1371/journal.pone.0256990.t001>

potential specifically for protein design and applying it to larger proteins is a topic of future work.

We also investigate whether the potential is able to distinguish between native structures and close decoys. Models up to 12 Å from the native structure were obtained using 3DRobot [66], which generates diverse and well-packed decoys using fragment assembly and energy minimisation. As shown in Fig 5B, the native structure has low energy compared to the decoys for Trp-cage and villin HP36. Chignolin was too small to run 3DRobot on. For BBA many decoys are lower in energy than the native structure. A number of these lower energy decoys are topological mirrors in which the  $\beta$  turn faces the other way, but the structure of the protein is otherwise native-like. The problem of mirror topologies has been discussed previously for protein structure determination [67] and de novo protein structure prediction [68].

## Discussion

The purpose of learning a coarse-grained force field for proteins is to demonstrate that DMS can be used to learn all the parameters from scratch in simple, interpretable force fields. It is notable that running training simulations from the native structure, reaching up to 4 Å RMSD to the native structure over 2,000 steps during training, is sufficient to learn a potential that can fold proteins from an extended chain of secondary structure elements over a few million steps. Whilst the nature of the training may give some bias to globular structures, as is the case with most force fields for proteins, we do notice that simulations at higher temperatures lead to unfolded and variable structures. The same potential can also be used for model scoring, despite the energy not being explicitly used at all beyond force calculation during training. This particular force field may be useful for exploring the conformational space of proteins, discovering folding pathways, assessing flexible or disordered regions of proteins, or predicting structure in combination with co-evolutionary or experimental constraints [69]. Other coarse-grained systems may be immediately amenable to DMS, for example protein aggregation [70]. The real possibility for DMS though lies in applying it to all-atom potentials with a variety of loss functions.

Here we have used RMSD as a simple loss function, but there are a variety of possible loss functions for DMS suitable for other systems. Examples include RMSD over the course of a simulation, radius of gyration, the radial distribution function, the flexibility of a set of atoms during the simulation, the distance between two molecules, supramolecular geometry (e.g. assembly of molecules into fibres), the correlation of different particle velocities, the energy of

a system, the temperature of a system, a measure of phase change, steered molecular dynamics, or some combination of the above. Many of the possible properties are based on static reference structures, thermodynamic observables or chemical knowledge, meaning that expensive trajectory data is not necessarily required. Complex constraints that might be difficult to use in the simulation itself can be targeted via the loss function. Any property that can be computed from the system with a meaningful gradient can be used to guide a learned force field to reproduce desired behaviour. This sets DMS apart from contrastive divergence approaches that train entire force fields by teaching them to distinguish native and non-native ensembles. Possible applications include improving force field accuracy on disordered proteins [8], combining DMS with learned atom typing [25] to study protein-ligand binding, improving torsion angle potentials to balance bonded and non-bonded terms [71], exploring whether multi-body terms can make molecular mechanics potentials more accurate, and developing potentials to promote protein-protein docking. The development of software packages appropriate for DMS such as JAX MD [72], TorchMD [18], DeePMD-kit [73], SchNetPack [74], and DiffTachi [75], along with the effort to make programming languages such as Julia differentiable by default [76], will assist in the development of DMS.

For DMS to be used to parameterise all-atom molecular mechanics force fields—either from scratch, by tuning existing parameters [8], or by adding new atom types—a few issues need to be addressed. Algorithms such as particle mesh Ewald for long-range electrostatic interactions [77] will have to be implemented in differentiable frameworks. The best form for temperature control during training will have to be considered, as stochasticity will likely affect the gradients. In this study, the thermostat was not used during training. Even the form of the numerical integration will have to be explored, as it is unclear that velocity Verlet integration with a standard time step is best-suited to DMS. In previous work damping of velocities or gradients have been used to prevent exploding gradients [34, 75]. It has also been shown that discrete time steps [75] and large time steps [34] can lead to incorrect gradients. Ensuring continuous potentials and reasonable time steps should prevent these issues, and the standard all-atom potentials—harmonic, cosine, Coulomb and Lennard-Jones—are all continuous. However, cutoffs for short-range forces and neighbour lists will have to be carefully considered. Finding the parameters for a simple water model would be an ideal system to start with [18].

DMS, like deep learning, appears to be sensitive to hyperparameters, and generally appropriate choices for these will have to be discovered. Learned potentials being amenable to physical interpretation helps when investigating such issues, as well as alleviating problems of under-specification and shortcut learning identified for neural networks [78]. For example, increasing the learning rate in this study gives jagged potentials, whereas decreasing it leads to prohibitively slow training.

As with analogous developments in deep learning, a limitation of DMS is the GPU memory required. Most deep learning software frameworks are geared towards reverse-mode AD, in which intermediate results of the forward pass of the network are stored and used during the reverse pass to calculate gradients. The requirement to store intermediate results means that it scales linearly in memory with the number of steps for DMS. By contrast, forward-mode AD does not store intermediate results because the gradients are calculated in tandem with the forward pass. The memory required does not therefore increase with the number of steps, though the number of learned parameters does affect the computation speed. Use of forward-mode AD may provide a way for DMS to use the large number of steps during training that would be required to learn molecular mechanics force fields. Advances in hardware, GPU parallelism and algorithmic techniques such as gradient checkpointing, invertible simulations [79], adjoint sensitivity [29, 30] and offloading compute to the CPU [80] also present solutions to the issue of GPU memory.

Though applications of deep learning in biology have been impressive, the models used have been largely taken off the shelf from other fields. DMS provides a general approach to improving and expanding the force fields that have been crucial for biological understanding. Previous limitations are rapidly being addressed by improvements in hardware, specialist software and the ever-increasing amount of experimental data available. This study shows that DMS can use the techniques and frameworks of neural network training but rely on proven, interpretable functional forms rather than deep neural networks themselves.

## Materials and methods

### Dataset

The PDB [61] was searched for protein chains with 20 to 100 residues, no internal missing residues and resolution 2.5 Å or better. These chains were clustered at 30% sequence identity to reduce redundancy. Proteins homologous to those used in the results were removed from the dataset by eliminating overlap at the same ECOD T-level [81] and removing hits when searching the PDB using BLAST with an E-value of 1.0. The resulting chains were randomly split into a training set of 2,004 chains and a validation set of 200 chains used to monitor training. Single sequence secondary structure prediction was carried out using PSIPRED [82]. Details of the proteins used in the results are given in S1 Table in [S1 File](#). Dataset collection and other aspects of the work were carried out using BioStructures.jl [83] and the Bio.PDB module of Biopython [84].

### Molecular simulation

The system was implemented in PyTorch [57]. A protein is represented in a coarse-grained manner with 4 point particles for each residue corresponding to backbone nitrogen, backbone C $\alpha$ , backbone carbonyl carbon and the centroid of the sidechain heavy atoms. There is no explicit solvent, no periodic boundary conditions and no neighbour list. The masses are set to 15 for N (includes amide H), 13 for C $\alpha$  (includes H) and 28 for C (includes carbonyl O). The sidechain mass for each amino acid is the sum of the atom masses in the all-atom sidechain, with glycine set artificially heavier at a mass of 10.

The overall learned potential used to obtain the forces at each time step consists of 3 components, each consisting of many individual potentials. Each individual potential consists of a number of energy values corresponding to 140 bins evenly distributed over a specified distance or angle range. Overall there are 29,360 individual potentials and 4,111,000 learnable parameters.

1. *Pairwise distance potentials.* There are 80 atom types (4 atoms for 20 amino acids) and each pair has a distance potential. In addition there are separate potentials for each atom pair on residues close in sequence with residue separations  $i \rightarrow i$  to  $i \rightarrow i + 4$ , allowing the model to learn local constraints separately from global preferences. Covalent bond interactions are included implicitly in the potentials for atom pairs on the same residue, and the  $i \rightarrow i + 1$  potentials in the case of the C-N backbone covalent bond. There are 28,960 pairwise distance potentials in total (3,240 general, 120 same residue,  $6,400 \times 4$  close residues). Each distance potential has 140 bins distributed between 1 Å and 15 Å (0.1 Å width per bin). For the close residue potentials the bins are distributed between 0.7 Å and 14.7 Å (0.1 Å width per bin), and for same residue pairs between 0.7 Å and 5.6 Å (0.035 Å width per bin).
2. *Bond angle potentials.* There are 5 bond angles in the model—3 in the backbone, 2 to the sidechain centroid—and there is a separate potential for each of these for each amino acid.

There are 100 bond angle potentials in total ( $5 \times 20$ ). Each angle potential has 140 bins distributed between  $60^\circ$  and  $180^\circ$  ( $0.86^\circ$  width per bin).

3. *Torsion angle potentials.* There are 5 torsion angles in the model—3 in the backbone, 2 to the sidechain centroid—and there is a separate potential for each of these for each amino acid. There are also separate potentials for residues predicted as  $\alpha$ -helical,  $\beta$ -sheet or coiled to help in secondary structure formation. There are 300 torsion angle potentials in total ( $5 \times 20 \times 3$ ). Each torsion angle potential has 140 bins distributed between  $-180^\circ$  and  $180^\circ$  ( $2.57^\circ$  width per bin) with an extra bin on either end to allow the derived force to be periodic.

At each simulation time step the force is calculated as the negative gradient of the potential using the following finite differencing procedure for each individual potential:

1. Calculate the current value of the property, e.g. the distance between two atoms or a bond angle.
2. Find the bin  $b_i$  with the closest bin centre to the value, excluding the first and last bins. For distances this means that all distances over the maximum bin distance (e.g.  $15 \text{ \AA}$ ) are treated as being in the penultimate bin.
3. Calculate  $F = \frac{1}{2}(E(b_{i-1}) - E(b_{i+1}))$ , where  $E(b_i)$  is the energy of bin  $b_i$ .
4. Multiply  $F$  by the appropriate vector on each atom to apply the force [85].

This is shown in Fig 1C. The sum of the resulting forces on each atom is divided by the atomic masses to get the accelerations. This approach is differentiable since it can be implemented in a vectorised manner in PyTorch. Whilst more sophisticated methods such as a sum of Gaussians or spline fitting could be adopted to obtain the force from the potential, this finite differencing procedure was found to be memory-efficient and effective. It was also found to be more effective than learning force values directly, an approach that does not immediately provide an interpretable potential that can be used to calculate energies.

The coordinates and velocities are updated at every time step using the velocity Verlet integrator. If  $\mathbf{x}(t)$  is the coordinates at time  $t$ ,  $\mathbf{v}(t)$  is the velocities at time  $t$ ,  $\mathbf{a}(t)$  is the accelerations at time  $t$ , and  $\Delta t$  is the time step, then the procedure is:

1. Calculate  $\mathbf{x}(t + \Delta t) = \mathbf{x}(t) + \mathbf{v}(t)\Delta t + \frac{1}{2}\mathbf{a}(t)\Delta t^2$ .
2. Obtain  $\mathbf{a}(t + \Delta t)$  from the learned potential as described above using  $\mathbf{x}(t + \Delta t)$ .
3. Calculate  $\mathbf{v}(t + \Delta t) = \mathbf{v}(t) + \frac{1}{2}(\mathbf{a}(t) + \mathbf{a}(t + \Delta t))\Delta t$ .
4. Update  $t + \Delta t$  to  $t$  and go to step 1.

No thermostat was used during training as it was not found to improve performance. Training therefore takes place in the NVE ensemble, with the caveat that the coarse-grained nature of the system means it does not have a conventional volume. A time step of 0.02 was used for training.

During production runs used to obtain results, which consisted of many more steps than during training, the Andersen thermostat was used to keep temperature constant [86]. Production runs therefore take place in the NVT ensemble. At each step, each atom is given a new velocity with probability equal to the time step divided by a coupling constant. The new velocities are drawn from a normal distribution with mean 0 and standard deviation equal to a temperature value. Simulations were run with temperature 0.015 and coupling constant 25, apart from the  $(\text{AAQAA})_3$  repeat peptide helix formation simulations which were run with

temperature 0.022 and coupling constant 100. A lower time step of 0.004 was used during production runs to ensure stability of the simulations. We did try using Langevin dynamics for production runs but found it did not improve the results.

## Training

During training, simulations are started from the native structure. Starting velocities are drawn from a normal distribution with mean 0 and standard deviation 0.1. At the end of each simulation  $R_f$  is calculated as the RMSD across all atoms in the coarse-grained model of the final conformation compared to the native structure using the Kabsch algorithm.  $\log(1 + R_f)$  was used as the loss function, which was found to give better performance than using  $R_f$ . PyTorch records operations during the simulation (the forward pass) in a directed acyclic graph. Once the loss is calculated, the graph is traversed backwards (the backward pass) and the known gradient functions at each step combined using the chain rule to obtain the gradients for each parameter. These gradients are used by the optimiser to modify the potential to reduce the  $R_f$  from future simulations. At the start of training all values in the potential are set to zero. We did try starting the potentials from PMFs derived from the PDB but this did not improve results. The Adam optimiser [87] was used with a learning rate of  $10^{-4}$ . Adam maintains different learning rates for each parameter and gave better results than stochastic gradient descent. Gradients are accumulated for 100 proteins before the optimiser updates the gradients. At epoch 38 the Adam optimiser was reset with a lower learning rate of  $5 \times 10^{-5}$ .

One epoch of training consists of simulating each protein in the training set in a random order. A batch size of 1 was used due to memory constraints. The number of steps in the simulation was increased over the course of training, starting at 250 for the first epoch and increasing by 250 every 5 epochs to a maximum of 2,000 at epoch 36 and beyond. This allowed the model to access lower  $R_f$  values during early epochs to learn basic chemical principles such as steric clashing and covalent bonding. Further increasing the step number was prevented by memory limitations of the GPU (32 GB); during training the memory required scales linearly with the number of steps due to the requirements of storing intermediate computations for reverse mode AD. This limitation is not present for production runs, which use memory constant in the number of steps and typically less than 1 GB. Training was carried out on a NVIDIA Tesla V100 for 45 epochs, which took around 2 months. This can likely be sped up using multiple GPUs. The training and validation  $R_f$  values throughout training are shown in S1 Fig in S1 File. The protein folding simulations shown in Fig 4B took around 36 hours for 12m steps, equating to around 10 ms per time step. This was constant across protein sizes tested due to the vectorised operations used. Simulation time is around 3 times slower on the CPU depending on the hardware used.

## Analysis

Throughout this study we analyse only coarse-grained models, however all-atom models can be generated using software such as PULCHRA [88] if required. Alanine dipeptide conformations were generated using PeptideBuilder [89] at  $\phi/\psi$  intervals of  $2^\circ$ . The (AAQAA)<sub>3</sub> repeat peptide simulation was started from a random conformation where each residue was given a  $\phi$  angle between  $-180^\circ$  and  $-30^\circ$  and a  $\psi$  angle between  $-180^\circ$  and  $180^\circ$ .  $\alpha$ -helical fraction was determined by the fraction of non-terminal residues where the  $\phi$  angle was between  $-120^\circ$  and  $-30^\circ$  and the  $\psi$  angle was between  $-60^\circ$  and  $30^\circ$ . The  $\alpha$ -helical fraction was averaged over a window of 5 snapshots either side of the snapshot in question for ease of visualisation. For simulations starting from predicted secondary structure, a residue was given an  $\alpha$ -helical starting

conformation ( $\phi -60^\circ$ ,  $\psi -60^\circ$ ) if predicted  $\alpha$ -helical and an extended starting conformation ( $\phi -120^\circ$ ,  $\psi 140^\circ$ ) if predicted  $\beta$ -sheet or coiled.

The representative structures in Fig 4B were found using MDAnalysis [90] and were visualised in PyMOL [91] after being run through PULCHRA. UNRES models were generated on the web server [63] using the parameters for the MREMD structure prediction example in the tutorial: the FF2 force field, extended chain start, secondary structure restraints, Berendsen thermostat with 1.0 coupling to the thermal bath, and 8 replicas exchanging every 1,000 steps with temperatures ranging from 270 K to 340 K in steps of 10 K. The number of steps was increased to the maximum of  $10^7$ . The top ranked model was used. CABS-fold models were generated on the web server [64] with default de novo parameters including CABS temperature 3.5–1.0. The top ranked model was used. In both cases we use the same single sequence secondary structure prediction as for our method.

The fraction of native sequence results involved threading a sequence onto the native structure. Sequences were chosen by mutating a given fraction of residues to amino acids taken from the background distribution of amino acids in the PDB. 30 sequences are generated for each fraction. Each sidechain centroid is placed at a distance from the  $C\alpha$  atom corresponding to the minimum in the learned potential for that amino acid, along the vector linking the  $C\alpha$  atom and the native sidechain centroid. A brief energy minimisation of 100 steps in the learned potential is carried out and the final energy is used as the energy of the sequence. For the fixed backbone design task 2,000 trial mutations were made. Each trial involved mutating one residue to an amino acid taken from the PDB distribution, calculating the new energy by threading the sequence and running energy minimisation as above, and accepting or rejecting the mutation. Mutations resulting in lower energy are always accepted, and mutations resulting in a higher energy of 10 or more energy units are always rejected. Mutations resulting in an energy increase up to 10 energy units are accepted with a probability that changes linearly from 0.25 at the start of the simulation to 0.0 at 1,000 trials, and remains at 0.0 for the remaining trials.

EvoEF2 [65] fixed backbone design runs were carried out with default de novo design parameters by running ‘EvoEF2 –command = ProteinDesign –monomer –pdb = input.pdb’. 3DRobot models were generated using the 3DRobot standalone software with default parameters [66]. Plots throughout were produced with Matplotlib [92], seaborn [93] and PyEMMA [94] for Fig 3B.

## Supporting information

**S1 File. Contains all the supporting tables and figures.**  
(PDF)

## Acknowledgments

We thank Nick Charron and Cecilia Clementi for providing the all-atom alanine dipeptide data from Wang et al. 2019 [15]. We thank the UCL Bioinformatics Group for useful discussions.

## Author Contributions

**Conceptualization:** Joe G. Greener.

**Data curation:** Joe G. Greener.

**Funding acquisition:** David T. Jones.

**Investigation:** Joe G. Greener.

**Methodology:** Joe G. Greener.

**Project administration:** David T. Jones.

**Software:** Joe G. Greener.

**Supervision:** David T. Jones.

**Validation:** Joe G. Greener.

**Writing – original draft:** Joe G. Greener.

**Writing – review & editing:** Joe G. Greener.

## References

1. Hollingsworth SA, Dror RO. Molecular Dynamics Simulation for All. *Neuron*. 2018; 99(6):1129–1143. <https://doi.org/10.1016/j.neuron.2018.08.011> PMID: 30236283
2. Brini E, Simmerling C, Dill K. Protein storytelling through physics. *Science*. 2020; 370:eaaz3041. <https://doi.org/10.1126/science.aaz3041> PMID: 33243857
3. Henzler-Wildman K, Kern D. Dynamic personalities of proteins. *Nature*. 2007; 450(7172):964–972. <https://doi.org/10.1038/nature06522> PMID: 18075575
4. Kmiecik S, Gront D, Kolinski M, Wieteska L, Dawid AE, Kolinski A. Coarse-Grained Protein Models and Their Applications. *Chem Rev*. 2016; 116(14):7898–7936. <https://doi.org/10.1021/acs.chemrev.6b00163> PMID: 27333362
5. Piana S, Lindorff-Larsen K, Shaw DE. How robust are protein folding simulations with respect to force field parameterization? *Biophys J*. 2011; 100(9):L47–L49. <https://doi.org/10.1016/j.bpj.2011.03.051> PMID: 21539772
6. Wang LP, Martinez TJ, Pande VS. Building Force Fields: An Automatic, Systematic, and Reproducible Approach. *J Phys Chem Lett*. 2014; 5(11):1885–1891. <https://doi.org/10.1021/jz500737m> PMID: 26273869
7. Best RB, Zheng W, Mittal J. Balanced Protein-Water Interactions Improve Properties of Disordered Proteins and Non-Specific Protein Association. *J Chem Theory Comput*. 2014; 10(11):5113–5124. <https://doi.org/10.1021/ct500569b> PMID: 25400522
8. Robustelli P, Piana S, Shaw DE. Developing a molecular dynamics force field for both folded and disordered protein states. *Proc Natl Acad Sci USA*. 2018; 115(21):E4758–E4766. <https://doi.org/10.1073/pnas.1800690115> PMID: 29735687
9. Liu M, Das AK, Lincoff J, Sasmal S, Cheng SY, Vernon R, et al. Configurational Entropy of Folded Proteins and its Importance for Intrinsically Disordered Proteins. *arXiv*. 2020;2007.06150.
10. Senior AW, Evans R, Jumper J, Kirkpatrick J, Sifre L, Green T, et al. Improved protein structure prediction using potentials from deep learning. *Nature*. 2020; 577(7792):706–710. <https://doi.org/10.1038/s41586-019-1923-7> PMID: 31942072
11. Noé F, De Fabritiis G, Clementi C. Machine learning for protein folding and dynamics. *Curr Opin Struct Biol*. 2020; 60:77–84. <https://doi.org/10.1016/j.sbi.2019.12.005> PMID: 31881449
12. Noé F, Tkatchenko A, Müller KR, Clementi C. Machine Learning for Molecular Simulation. *Annu Rev Phys Chem*. 2020; 71:361–390. <https://doi.org/10.1146/annurev-physchem-042018-052331> PMID: 32092281
13. Gkeka P, Stoltz G, Barati Farimani A, Belkacemi Z, Ceriotti M, Chodera JD, et al. Machine Learning Force Fields and Coarse-Grained Variables in Molecular Dynamics: Application to Materials and Biological Systems. *J Chem Theory Comput*. 2020; 16(8):4757–4775. <https://doi.org/10.1021/acs.jctc.0c00355> PMID: 32559068
14. Zhang L, Han J, Wang H, Car R, E W. DeePCG: Constructing coarse-grained models via deep neural networks. *J Chem Phys*. 2018; 149(3):034101. <https://doi.org/10.1063/1.5027645> PMID: 30037247
15. Wang J, Olsson S, Wehmeyer C, Pérez A, Charron NE, De Fabritiis G, et al. Machine Learning of Coarse-Grained Molecular Dynamics Force Fields. *ACS Cent Sci*. 2019; 5(5):755–767. <https://doi.org/10.1021/acscentsci.8b00913> PMID: 31139712
16. Husic BE, Charron NE, Lemm D, Wang J, Pérez A, Majewski M, et al. Coarse graining molecular dynamics with graph neural networks. *J Chem Phys*. 2020; 153(19):194101. <https://doi.org/10.1063/5.0026133> PMID: 33218238



17. Wang W, Gómez-Bombarelli R. Coarse-graining auto-encoders for molecular dynamics. *npj Computational Materials*. 2019; 5(125).
18. Doerr S, Majewski M, Pérez A, Krämer A, Clementi C, Noé F, et al. TorchMD: A deep learning framework for molecular simulations. *arXiv*. 2020;2012.12106.
19. Chmiela S, Tkatchenko A, Sauceda HE, Poltavsky I, Schütt KT, Müller KR. Machine learning of accurate energy-conserving molecular force fields. *Sci Adv*. 2017; 3(5):e1603015. <https://doi.org/10.1126/sciadv.1603015> PMID: 28508076
20. Bogojeski M, Vogt-Maranto L, Tuckerman ME, Müller KR, Burke K. Quantum chemical accuracy from density functional approximations via machine learning. *Nat Commun*. 2020; 11:5223. <https://doi.org/10.1038/s41467-020-19093-1> PMID: 33067479
21. Hermann J, Schätzle Z, Noé F. Deep-neural-network solution of the electronic Schrödinger equation. *Nat Chem*. 2020; 12(10):891–897. <https://doi.org/10.1038/s41557-020-0544-y> PMID: 32968231
22. Batzner S, Smidt TE, Sun L, Mailoa JP, Kornbluth M, Molinari N, et al. SE(3)-Equivariant Graph Neural Networks for Data-Efficient and Accurate Interatomic Potentials. *arXiv*. 2021;2101.03164.
23. Noé F, Olsson S, Köhler J, Wu H. Boltzmann generators: Sampling equilibrium states of many-body systems with deep learning. *Science*. 2019; 365:eaaw1147. <https://doi.org/10.1126/science.aaw1147> PMID: 31488660
24. Shin K, Tran DP, Takemura K, Kitao A, Terayama K, Tsuda K. Enhancing Biomolecular Sampling with Reinforcement Learning: A Tree Search Molecular Dynamics Simulation Method. *ACS Omega*. 2019; 4(9):13853–13862. <https://doi.org/10.1021/acsomega.9b01480> PMID: 31497702
25. Wang Y, Fass J, Chodera JD. End-to-End Differentiable Molecular Mechanics Force Field Construction. *arXiv*. 2020;2010.01196.
26. AlQuraishi M. End-to-End Differentiable Learning of Protein Structure. *Cell Systems*. 2019; 8(4):292–301.e3. <https://doi.org/10.1016/j.cels.2019.03.006> PMID: 31005579
27. Ragoza M, Hochuli J, Idrobo E, Sunseri J, Koes DR. Protein-Ligand Scoring with Convolutional Neural Networks. *J Chem Inf Model*. 2017; 57(4):942–957. <https://doi.org/10.1021/acs.jcim.6b00740> PMID: 28368587
28. Rufa DA, Bruce Macdonald HE, Fass J, Wieder M, Grinaway PB, Roitberg AE, et al. Towards chemical accuracy for alchemical free energy calculations with hybrid physics-based machine learning / molecular mechanics potentials. *bioRxiv*. 2020;<https://www.biorxiv.org/content/10.1101/2020.07.29.227959v1>.
29. Wang W, Axelrod S, Gómez-Bombarelli R. Differentiable Molecular Simulations for Control and Learning. *arXiv*. 2020;2003.00868.
30. Chen RTQ, Rubanova Y, Bettencourt J, Duvenaud D. Neural Ordinary Differential Equations. *NeurIPS*. 2018;<https://arxiv.org/abs/1806.07366>.
31. Rackauckas C, Ma Y, Martensen J, Warner C, Zubov K, Supekar R, et al. Universal Differential Equations for Scientific Machine Learning. *arXiv*. 2020;2001.04385.
32. Holl P, Koltun V, Thuerey N. Learning to Control PDEs with Differentiable Physics. *arXiv*. 2020;2001.07457.
33. Lindorff-Larsen K, Maragakis P, Piana S, Eastwood MP, Dror RO, Shaw DE. Systematic validation of protein force fields against experimental data. *PLoS ONE*. 2012; 7(2):e32131. <https://doi.org/10.1371/journal.pone.0032131> PMID: 22384157
34. Ingraham J, Riesselman A, Sander C, Marks D. Learning Protein Structure with a Differentiable Simulator. *ICLR*. 2019;<https://openreview.net/forum?id=Byg3y3C9Km>.
35. Jumper JM, Faruk NF, Freed KF, Sosnick TR. Trajectory-based training enables protein simulations with accurate folding and Boltzmann ensembles in cpu-hours. *PLoS Comput Biol*. 2018; 14(12):e1006578. <https://doi.org/10.1371/journal.pcbi.1006578> PMID: 30589834
36. Várnai C, Burkoff NS, Wild DL. Efficient Parameter Estimation of Generalizable Coarse-Grained Protein Force Fields Using Contrastive Divergence: A Maximum Likelihood Approach. *J Chem Theory Comput*. 2013; 9(12):5718–5733. <https://doi.org/10.1021/ct400628h> PMID: 24683370
37. Podtelezhnikov AA, Ghahramani Z, Wild DL. Learning about protein hydrogen bonding by minimizing contrastive divergence. *Proteins*. 2007; 66(3):588–599. <https://doi.org/10.1002/prot.21247> PMID: 17109405
38. Zaborowski B, Jagieła D, Czaplewski C, Hałabis A, Lewandowska A, Zmudzińska W, et al. A Maximum-Likelihood Approach to Force-Field Calibration. *J Chem Inf Model*. 2015; 55(9):2050–2070. <https://doi.org/10.1021/acs.jcim.5b00395> PMID: 26263302
39. Krupa P, Hałabis A, Zmudzińska W, Ołdziej S, Scheraga HA, Liwo A. Maximum Likelihood Calibration of the UNRES Force Field for Simulation of Protein Structure and Dynamics. *J Chem Inf Model*. 2017; 57(9):2364–2377. <https://doi.org/10.1021/acs.jcim.7b00254> PMID: 28809487

40. Winther O, Krogh A. Teaching computers to fold proteins. *Physical Review E*. 2004; 70(3):030903. <https://doi.org/10.1103/PhysRevE.70.030903> PMID: 15524499
41. Fain B, Levitt M. Funnel sculpting for in silico assembly of secondary structure elements of proteins. *Proc Natl Acad Sci USA*. 2003; 100(19):10700–10705. <https://doi.org/10.1073/pnas.1732312100> PMID: 12925740
42. Park H, Zhou G, Baek M, Baker D, DiMaio F. Force Field Optimization Guided by Small Molecule Crystal Lattice Data Enables Consistent Sub-Angstrom Protein-Ligand Docking. *J Chem Theory Comput*. 2021; 17:2000–2010. <https://doi.org/10.1021/acs.jctc.0c01184> PMID: 33577321
43. Crippen GM, Snow ME. A 1.8 Å resolution potential function for protein folding. *Biopolymers*. 1990; 29(10-11):1479–1489. <https://doi.org/10.1002/bip.360291014> PMID: 2361157
44. Fujitsuka Y, Takada S, Luthey-Schulten ZA, Wolynes PG. Optimizing physical energy functions for protein folding. *Proteins*. 2004; 54(1):88–103. <https://doi.org/10.1002/prot.10429> PMID: 14705026
45. Demerdash O, Shrestha UR, Petridis L, Smith JC, Mitchell JC, Ramanathan A. Using Small-Angle Scattering Data and Parametric Machine Learning to Optimize Force Field Parameters for Intrinsically Disordered Proteins. *Front Mol Biosci*. 2019; 6:64. <https://doi.org/10.3389/fmolb.2019.00064> PMID: 31475155
46. Baydin AG, Pearlmutter BA, Radul AA, Siskind JM. Automatic differentiation in machine learning: a survey. *Journal of Machine Learning Research*. 2018; 18(1):1–43.
47. Liwo A, Khalili M, Scheraga HA. Ab initio simulations of protein-folding pathways by molecular dynamics with the united-residue model of polypeptide chains. *Proc Natl Acad Sci USA*. 2005; 102(7):2362–2367. <https://doi.org/10.1073/pnas.0408885102> PMID: 15677316
48. Kolinski A. Protein modeling and structure prediction with a reduced representation. *Acta Biochim Pol*. 2004; 51(2):349–371. [https://doi.org/10.18388/abp.2004\\_3575](https://doi.org/10.18388/abp.2004_3575) PMID: 15218533
49. Hubner IA, Deeds EJ, Shakhnovich EI. High-resolution protein folding with a transferable potential. *Proc Natl Acad Sci USA*. 2005; 102(52):18914–18919. <https://doi.org/10.1073/pnas.0502181102> PMID: 16365306
50. Izvekov S, Voth GA. A multiscale coarse-graining method for biomolecular systems. *J Phys Chem B*. 2005; 109(7):2469–2473. <https://doi.org/10.1021/jp044629q> PMID: 16851243
51. Maupetit J, Tuffery P, Derreumaux P. A coarse-grained protein force field for folding and structure prediction. *Proteins*. 2007; 69(2):394–408. <https://doi.org/10.1002/prot.21505> PMID: 17600832
52. Zhou H, Zhou Y. Distance-scaled, finite ideal-gas reference state improves structure-derived potentials of mean force for structure selection and stability prediction. *Protein Sci*. 2002; 11(11):2714–2726. <https://doi.org/10.1110/ps.0217002> PMID: 12381853
53. Shen MY, Sali A. Statistical potential for assessment and prediction of protein structures. *Protein Sci*. 2006; 15(11):2507–2524. <https://doi.org/10.1110/ps.062416606> PMID: 17075131
54. Innes M, Edelman A, Fischer K, Rackauckas C, Saba E, Shah VB, et al. A Differentiable Programming System to Bridge Machine Learning and Scientific Computing. *arXiv*. 2019;1907.07587.
55. Goodrich CP, King EM, Schoenholz SS, Cubuk ED, Brenner M. Self-assembling kinetics: Accessing a new design space via differentiable statistical-physics models. *arXiv*. 2020;2010.15175.
56. Li L, Hoyer S, Pederson R, Sun R, Cubuk ED, Riley P, et al. Kohn-Sham Equations as Regularizer: Building Prior Knowledge into Machine-Learned Physics. *Phys Rev Lett*. 2021; 126(3):036401. PMID: 33543980
57. Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, et al. PyTorch: An Imperative Style, High-Performance Deep Learning Library. *Advances in Neural Information Processing Systems* 32. 2019; p. 8024–8035.
58. Krieger E, Koraimann G, Vriend G. Increasing the precision of comparative models with YASARA NOVA—a self-parameterizing force field. *Proteins*. 2002; 47(3):393–402. <https://doi.org/10.1002/prot.10104> PMID: 11948792
59. Lindorff-Larsen K, Piana S, Dror RO, Shaw DE. How fast-folding proteins fold. *Science*. 2011; 334(6055):517–520. <https://doi.org/10.1126/science.1208351> PMID: 22034434
60. Nguyen H, Maier J, Huang H, Perrone V, Simmerling C. Folding simulations for proteins with diverse topologies are accessible in days with a physics-based force field and implicit solvent. *J Am Chem Soc*. 2014; 136(40):13959–13962. <https://doi.org/10.1021/ja5032776> PMID: 25255057
61. Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, et al. The Protein Data Bank. *Nucleic Acids Res*. 2000; 28(1):235–242. <https://doi.org/10.1093/nar/28.1.235> PMID: 10592235
62. Honda S, Akiba T, Kato YS, Sawada Y, Sekijima M, Ishimura M, et al. Crystal structure of a ten-amino acid protein. *J Am Chem Soc*. 2008; 130(46):15327–15331. <https://doi.org/10.1021/ja8030533> PMID: 18950166

63. Czaplewski C, Karczynska A, Sieradzan AK, Liwo A. UNRES server for physics-based coarse-grained simulations and prediction of protein structure, dynamics and thermodynamics. *Nucleic Acids Res.* 2018; 46(W1):W304–W309. <https://doi.org/10.1093/nar/gky328> PMID: 29718313
64. Blaszczyk M, Jamroz M, Kmiecik S, Kolinski A. CABS-fold: Server for the de novo and consensus-based prediction of protein structure. *Nucleic Acids Res.* 2013; 41(W1):W406–W411. <https://doi.org/10.1093/nar/gkt462> PMID: 23748950
65. Huang X, Pearce R, Zhang Y. EvoEF2: accurate and fast energy function for computational protein design. *Bioinformatics.* 2020; 36(4):1135–1142. <https://doi.org/10.1093/bioinformatics/btz740> PMID: 31588495
66. Deng H, Jia Y, Zhang Y. 3DRobot: automated generation of diverse and well-packed protein structure decoys. *Bioinformatics.* 2016; 32(3):378–387. <https://doi.org/10.1093/bioinformatics/btv601> PMID: 26471454
67. Pastore A, Atkinson RA, Saudek V, Williams RJ. Topological mirror images in protein structure computation: an underestimated problem. *Proteins.* 1991; 10(1):22–32. <https://doi.org/10.1002/prot.340100104> PMID: 1648217
68. Greener JG, Kandathil SM, Jones DT. Deep learning extends de novo protein modelling coverage of genomes using iteratively predicted structural constraints. *Nat Commun.* 2019; 10(1):3977. <https://doi.org/10.1038/s41467-019-11994-0> PMID: 31484923
69. Cheung NJ, Yu W. De novo protein structure prediction using ultra-fast molecular dynamics simulation. *PLoS ONE.* 2018; 13(11):e0205819. <https://doi.org/10.1371/journal.pone.0205819> PMID: 30458007
70. Šarić A, Chebaro YC, Knowles TP, Frenkel D. Crucial role of nonspecific interactions in amyloid nucleation. *Proc Natl Acad Sci USA.* 2014; 111(50):17869–17874. <https://doi.org/10.1073/pnas.1410159111> PMID: 25453085
71. Nerenberg PS, Head-Gordon T. New developments in force fields for biomolecular simulations. *Curr Opin Struct Biol.* 2018; 49:129–138. <https://doi.org/10.1016/j.sbi.2018.02.002> PMID: 29477047
72. Schoenholz SS, Cubuk ED. JAX, M.D.: A Framework for Differentiable Physics. *arXiv.* 2019;1912.04232.
73. Wang H, Zhang L, Han J, E W. DeePMD-kit: A deep learning package for many-body potential energy representation and molecular dynamics. *Computer Physics Communications.* 2018; 228:178–184. <https://doi.org/10.1016/j.cpc.2018.03.016>
74. Schütt KT, Kessel P, Gastegger M, Nicoli KA, Tkatchenko A, Müller KR. SchNetPack: A Deep Learning Toolbox For Atomistic Systems. *J Chem Theory Comput.* 2019; 15(1):448–455. <https://doi.org/10.1021/acs.jctc.8b00908> PMID: 30481453
75. Hu Y, Anderson L, Li TM, Sun Q, Carr N, Ragan-Kelley J, et al. DiffTaichi: Differentiable Programming for Physical Simulation. *arXiv.* 2019;1910.00935.
76. Innes M. Don't Unroll Adjoint: Differentiating SSA-Form Programs. *arXiv.* 2018;1810.07951.
77. Darden T, York D, Pedersen L. Particle mesh Ewald: An  $N \log(N)$  method for Ewald sums in large systems. *J Chem Phys.* 1993; 98(12):10089. <https://doi.org/10.1063/1.464397>
78. Geirhos R, Jacobsen JH, Michaelis C, Zemel R, Brendel W, Bethge M, et al. Shortcut learning in deep neural networks. *Nat Mach Intell.* 2020; 2:665–673. <https://doi.org/10.1038/s42256-020-00257-z>
79. Ardizzone L, Kruse J, Wirkert S, Rahner D, Pellegrini EW, Klessen RS, et al. Analyzing inverse problems with invertible neural networks. *ICLR.* 2019;<https://openreview.net/forum?id=rJed6jOckX>.
80. Ren J, Rajbhandari S, Aminabadi RY, Ruwase O, Yang S, Zhang M, et al. ZeRO-Offload: Democratizing Billion-Scale Model Training. *arXiv.* 2021;2101.06840.
81. Cheng H, Schaeffer RD, Liao Y, Kinch LN, Pei J, Shi S, et al. ECOD: an evolutionary classification of protein domains. *PLoS Comput Biol.* 2014; 10(12):e1003926. <https://doi.org/10.1371/journal.pcbi.1003926> PMID: 25474468
82. Jones DT. Protein secondary structure prediction based on position-specific scoring matrices. *J Mol Biol.* 1999; 292(2):195–202. <https://doi.org/10.1006/jmbi.1999.3091> PMID: 10493868
83. Greener JG, Selvaraj J, Ward BJ. BioStructures.jl: read, write and manipulate macromolecular structures in Julia. *Bioinformatics.* 2020; 36(14):4206–4207. <https://doi.org/10.1093/bioinformatics/btaa502> PMID: 32407511
84. Hamelryck T, Manderick B. PDB file parser and structure class implemented in Python. *Bioinformatics.* 2003; 19(17):2308–2310. <https://doi.org/10.1093/bioinformatics/btg299> PMID: 14630660
85. Monasse B, Boussinot F. Determination of Forces from a Potential in Molecular Dynamics. *arXiv.* 2014;1401.1181.
86. Andersen HC. Molecular dynamics simulations at constant pressure and/or temperature. *The Journal of Chemical Physics.* 1980; 72(4):2384–2393. <https://doi.org/10.1063/1.439486>

87. Kingma DP, Ba JL. Adam: A Method for Stochastic Optimization. ICLR. 2015;<https://arxiv.org/abs/1412.6980>.
88. Rotkiewicz P, Skolnick J. Fast procedure for reconstruction of full-atom protein models from reduced representations. *J Comput Chem*. 2008; 29(9):1460–1465. <https://doi.org/10.1002/jcc.20906> PMID: 18196502
89. Tien MZ, Sydykova DK, Meyer AG, Wilke CO. PeptideBuilder: A simple Python library to generate model peptides. *PeerJ*. 2013; 1:e80. <https://doi.org/10.7717/peerj.80> PMID: 23717802
90. Gowers RJ, Linke M, Barnoud J, Reddy TJE, Melo MN, Seyler SL, et al. MDAnalysis: A Python Package for the Rapid Analysis of Molecular Dynamics Simulations. *Proceedings of the 15th Python in Science Conference*. 2016; p. 98–105.
91. Schrödinger, LLC. The PyMOL Molecular Graphics System; 2020.
92. Hunter JD. Matplotlib: A 2D graphics environment. *Computing in Science & Engineering*. 2007; 9(3):90–95. <https://doi.org/10.1109/MCSE.2007.55>
93. Waskom M, the seaborn development team. mwaskom/seaborn. Zenodo. 2020;<https://doi.org/10.5281/zenodo.592845>.
94. Scherer MK, Trendelkamp-Schroer B, Paul F, Pérez-Hernández G, Hoffmann M, Plattner N, et al. PyEMMA 2: A Software Package for Estimation, Validation, and Analysis of Markov Models. *J Chem Theory Comput*. 2015; 11(11):5525–5542. <https://doi.org/10.1021/acs.jctc.5b00743> PMID: 26574340