



# Review Deep Learning in Protein Structural Modeling and Design

Wenhao Gao,<sup>1</sup> Sai Pooja Mahajan,<sup>1</sup> Jeremias Sulam,<sup>2</sup> and Jeffrey J. Gray<sup>1,\*</sup>

<sup>1</sup>Department of Chemical and Biomolecular Engineering, Johns Hopkins University, Baltimore, MD 21218, USA <sup>2</sup>Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD 21218, USA

https://doi.org/10.1016/j.patter.2020.100142

**THE BIGGER PICTURE** Proteins are linear polymers that fold into an incredible variety of three-dimensional structures that enable sophisticated functionality for biology. Computational modeling allows scientists to predict the three-dimensional structure of proteins from genomes, predict properties or behavior of a protein, and even modify or design new proteins for a desired function. Advances in machine learning, especially deep learning, are catalyzing a revolution in the paradigm of scientific research. In this review, we summarize recent work in applying deep learning techniques to tackle problems in protein structural modeling and design. Some deep learning-based approaches, especially in structure prediction, now outperform conventional methods, often in combination with higher-resolution physical modeling. Challenges remain in experimental validation, benchmarking, leveraging known physics and interpreting models, and extending to other biomolecules and contexts.

**345 Development/Pre-production:** Data science output has been rolled out/validated across multiple domains/problems

#### SUMMARY

1

2

Deep learning is catalyzing a scientific revolution fueled by big data, accessible toolkits, and powerful computational resources, impacting many fields, including protein structural modeling. Protein structural modeling, such as predicting structure from amino acid sequence and evolutionary information, designing proteins toward desirable functionality, or predicting properties or behavior of a protein, is critical to understand and engineer biological systems at the molecular level. In this review, we summarize the recent advances in applying deep learning techniques to tackle problems in protein structural modeling and design. We dissect the emerging approaches using deep learning techniques for protein structural modeling and discuss advances and challenges that must be addressed. We argue for the central importance of structure, following the "sequence  $\rightarrow$  structure  $\rightarrow$  function" paradigm. This review is directed to help both computational biologists to gain familiarity with the deep learning methods applied in protein modeling, and computer scientists to gain perspective on the biologically meaningful problems that may benefit from deep learning techniques.

#### **INTRODUCTION**

Proteins are linear polymers that fold into various specific conformations to function. The incredible variety of three-dimensional (3D) structures determined by the combination and order in which 20 amino acids thread the protein polymer chain (sequence of the protein) enables the sophisticated functionality of proteins responsible for most biological activities. Hence, obtaining the structures of proteins is of paramount importance in both understanding the fundamental biology of health and disease and developing therapeutic molecules. While protein structure is primarily determined by sophisticated experimental techniques, such as X-ray crystallography,<sup>1</sup> NMR spectroscopy<sup>2</sup> and, increasingly, cryoelectron microscopy,<sup>3</sup> computational structure prediction from the genetically encoded amino acid sequence of a protein has been used as an alternative when experimental approaches are limited. Computational methods have been used to predict the structure of proteins,<sup>4</sup> illustrate the mechanism of biological processes,<sup>5</sup> and determine the properties of proteins.<sup>6</sup> Furthermore, all naturally occurring proteins are a result of an evolutionary process of random variants arising under various selective pressures. Through this process, nature has explored only a small subset of theoretically possible protein sequence space. To explore a broader sequence and

Check fo

<sup>\*</sup>Correspondence: jgray@jhu.edu





#### Figure 1. Striking Improvement in Model Accuracy in CASP13 Due to the Deployment of Deep Learning Methods

(A) Trend lines of backbone accuracy for the best models in each of the 13 CASP experiments. Individual target points are shown for the two most recent experiments. The accuracy metric, GDT\_TS, is a multiscale indicator of the closeness of the C $\alpha$  atoms in a model to those in the corresponding experimental structure (higher numbers are more accurate). Target difficulty is based on sequence and structure similarity to other proteins with known experimental structures (see Kryshtafovych et al.<sup>4</sup> for details). Figure from Kryshtafovych et al. (2019).<sup>4</sup>

(B) Number of FM + FM/TBM (FM, free modeling; TBM, template-based modeling) domains (out of 43) solved to a TM score threshold for all groups in CASP.<sup>13</sup> AlphaFold ranked first among them, showing that the progress is mainly due to the development of DL-based methods. Figure from Senior et al. (2020).<sup>26</sup>

structural space that potentially contains proteins with enhanced or novel properties, techniques, such as *de novo* design can be used to generate new biological molecules that have the potential to tackle many outstanding challenges in biomedicine and biotechnology.<sup>7,8</sup>

While the application of machine learning and more general statistical methods in protein modeling can be traced back decades,<sup>9–13</sup> recent advances in machine learning, especially in deep learning (DL)-related techniques,<sup>14</sup> have opened up new avenues in many areas of protein modeling.<sup>15–18</sup> DL is a set of machine learning techniques based on stacked neural network layers that parameterize functions in terms of compositions of affine transformations and non-linear activation functions. Their ability to extract domain-specific features that are adaptively learned from data for a particular task often enables them to surpass the performance of more traditional methods. DL has made dramatic impacts on digital applications like image classification,<sup>19</sup> speech recognition,<sup>20</sup> and game playing.<sup>21</sup> Success in these areas has inspired an increasing interest in more complex data types, including protein structures.<sup>22</sup> In the most recent Critical Assessment of Structure Prediction (CASP13 held in 2018),<sup>4</sup> a biennial community experiment to determine the state-of-the-art in protein structure prediction, DL-based methods accomplished a striking improvement in model accuracy (see Figure 1), especially in the "difficult" target category where comparative modeling (starting with a known, related structure) is ineffective. The CASP13 results show that the complex mapping from amino acid sequence to 3D protein structure can be successfully learned by a neural network and generalized to unseen cases. Concurrently, for the protein design problem, progress in the field of deep generative models has spawned a range of promising approaches.<sup>23-25</sup>

In this review, we summarize the recent progress in applying DL techniques to the problem of protein modeling and discuss the potential pros and cons. We limit our scope to protein structure and function prediction, protein design with DL (see Figure 2), and a wide array of popular frameworks used in these applications. We discuss the importance of protein representation, and summarize the approaches to protein design based on DL for the first time. We also emphasize the central importance of protein structure, following the sequence  $\rightarrow$  structure  $\rightarrow$  function paradigm and argue that approaches based on structures may be most fruitful. We refer the reader to other review papers for more information on applications of DL in biology and medicine,<sup>16,15</sup> bioinformatics,<sup>27</sup> structural biology,<sup>17</sup> folding and dynamics, <sup>18,28</sup> antibody modeling,<sup>29</sup> and structural annotation and prediction of proteins.<sup>30,31</sup> Because DL is a fast-moving. interdisciplinary field, we chose to include preprints in this review. We caution the reader that these contributions have not been peer-reviewed, yet are still worthy of attention now for their ideas. In fact, in communities such as computer science, it is not uncommon for manuscripts to remain in this stage indefinitely. and some seminal contributions, such as Kingma and Welling's definitive paper on autoencoders (AEs),<sup>32</sup> are only available as preprints. In addition, we urge caution with any protein design studies that are purely in silico, and we highlight those that include experimental validation as a sign of their trustworthiness.

#### PROTEIN STRUCTURE PREDICTION AND DESIGN

#### **Problem Definition**

The prediction of protein 3D structure from amino acid sequence has been a grand challenge in computational biophysics for decades.<sup>33,34</sup> Folding of peptide chains is a fundamental concept in biophysics, and atomic-level structures of proteins and complexes are often the starting point to understand their function and to modulate or engineer them. Thanks to the recent advances in next-generation sequencing technology, there are now over 180 million protein sequences recorded in the UniProt dataset.<sup>35</sup> In contrast, only 158,000 experimentally determined structures are available in the Protein Data Bank. Thus, computational structure prediction is a critical problem of both practical and theoretical interest.

More recently, the advances in structure prediction have led to an increasing interest in the protein design problem. In design,

### Patterns Review



## CellPress OPEN ACCESS

Figure 2. Schematic Comparison of Three Major Tasks in Protein Modeling: Function Prediction, Structure Prediction, and Protein Design

In function prediction, the sequence and/or the structure is known and the functionality is needed as output of a neural net. In structure prediction, sequence is known input and structure is unknown output. Protein design starts from desired functionality, or a step further, structure that can perform this functionality. The desired output is a sequence that can fold into the structure or has such functionality.

the objective is to obtain a novel protein sequence that will fold into a desired structure or perform a specific function, such as catalysis. Naturally occurring proteins represent only an infinitesimal subset of all possible amino acid sequences selected by the evolutionary process to perform a specific biological function. Proteins with more robustness (higher thermal stability, resistance to degradation) or enhanced properties (faster catalysis, tighter binding) might lie in the space that has not been explored by nature, but is potentially accessible by de novo design. The current approach for computational de novo design is based on physical and evolutionary principles and requires significant domain expertise. Some successful examples include novel folds,<sup>36</sup> enzymes,<sup>37</sup> vaccines,<sup>38</sup> novel protein assemblies,<sup>39</sup> ligand-binding protein,<sup>40</sup> and membrane proteins.<sup>41</sup> While some papers occasionally refer to redesign of naturally occurring proteins or interfaces as "de novo", in this review we restrict that term only to works where completely new folds or interfaces are created.

#### **Conventional Computational Approaches**

The current methodology for computational protein structure prediction is largely based on Anfinsen's<sup>42</sup> thermodynamic hypothesis, which states that the native structure of a protein must be the one with the lowest free energy, governed by the energy landscape of all possible conformations associated with its sequence. Finding the lowest-energy state is challenging because of the immense space of possible conformations available to a protein, also known as the "sampling problem" or Levinthal's<sup>43</sup> paradox. Furthermore, the approach requires accurate free energy functions to describe the protein energy landscape and rank different conformations based on their energy, referred to as the "scoring problem." In light of these challenges, current computational techniques rely heavily on multiscale approaches. Low-resolution, coarse-grained energy functions are used to capture large-scale conformational sampling, such as the hydrophobic burial and formation of local secondary structural elements. Higher-resolution energy functions are used to explicitly model finer details, such as amino acid side-chain packing, hydrogen bonding, and salt bridges.<sup>44</sup>

Protein design problems, sometimes known as the inverse of structure prediction problems, require a similar toolbox. Instead of sampling the conformational space, a protein design protocol samples the sequence space that folds into the desired topology. Past efforts can be broadly divided into two broad classes: modifying an existing protein with known sequence and properties, or generating novel proteins with sequences and/or folds unrelated to those found in nature. The former class evolves an existing protein's amino acid sequence (and as a result, structure and properties) and can be loosely referred to as protein engineering or protein redesign. The latter class of methods is called de novo protein design, a term originally coined in 1997 when Dahiyat and Mayo<sup>45</sup> designed the FSD-1 protein, a soluble protein with a completely new sequence that folded into the previously known structure of a zinc finger. Korendovych and De-Grado's<sup>46</sup> recent retrospective chronicles the development of de novo design. Originally de novo design meant creation of entirely new proteins from scratch exploiting a target structure but, especially in the DL era, many authors now use the term to include methods that ignore structure in creating new sequences, often using extensive training data from known proteins in a particular functional class. In this review, we split our discussion of methods according to whether they trained directly between sequence and function (as certain natural language processing [NLP]-based DL paradigms allow), or whether they directly include protein structural data (like historical methods in rational protein design; see below in the section on "Protein Design").

Despite significant progress in the last several decades in the field of computational protein structure prediction and design,<sup>7,34</sup> accurate structure prediction and reliable design both remain challenging. Conventional approaches rely heavily on the accuracy of the energy functions to describe protein physics and the efficiency of sampling algorithms to explore the immense protein sequence and structure space. Both protein engineering and *de novo* approaches are often combined with experimental directed evolution<sup>8,47</sup> to achieve the optimal final molecules.<sup>7</sup>

#### **DL ARCHITECTURES**

In conventional computational approaches, predictions from data are made by means of physical equations and modeling. Machine learning puts forward a different paradigm in which algorithms automatically infer-or learn-a relationship between inputs and outputs from a set of hypotheses. Consider a collection of N training samples comprising features x in an input space  $\mathcal{X}$  (e.g., amino acid sequences), and corresponding labels y in some output space  $\mathcal{Y}$  (e.g., residue pairwise distances), where  $\{\mathbf{x}_i, y_i\}_{i=1}^N$  are sampled independently and identically distributed from some joint distribution  $\mathcal{P}$ . In addition, consider a function  $f: \mathcal{X} \rightarrow \mathcal{Y}$  in some function class  $\mathcal{H}$ , and a loss function  $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  that measures how much  $f(\mathbf{x})$  deviates from the corresponding label y. The goal of supervised learning is to find a function  $f \in \mathcal{H}$  that minimizes the expected loss,  $\mathbb{E}[\ell(f(\mathbf{x}))]$ y)], for  $(\mathbf{x}, \mathbf{y})$  sampled from  $\mathcal{P}$ . Since one does not have access to the true distribution but rather N samples from it, the







(A) CNNs are widely used in structure prediction.

(B) RNNs learn in an auto-regressive way and can be used for sequence generation.

(C) The VAE can be jointly trained by protein and properties to construct a latent space correlated with properties.

(D) In the GAN setting, a mapping from a priori distribution to the design space can be obtained via the adversarial training.

popular empirical risk minimization (ERM) approach seeks to minimize the loss over the training samples instead. In neural network models, in particular, the function class is parameterized by a collection of weights. Denoting these parameters collectively by  $\theta$ , ERM boils down to an optimization problem of the form

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^{N} \ell(f_{\theta}(\boldsymbol{x}_{i}), y_{i}).$$
 (Equation 1)

The choice of the network determines how the hypothesis class is parameterized. Deep neural networks typically implement a non-linear function as the composition of affine maps,  $W_l : \mathbb{R}^{n_l} \to \mathbb{R}^{n_{l+1}}$ , where  $W_l \mathbf{x} = W_l \mathbf{x} + \mathbf{b}_l$ , and other non-linear activation functions,  $\sigma(\cdot)$ . Rectifying linear units and max-pooling are some of the most popular non-linear transformations applied in practice. The architecture of the most popular option being their sequential composition  $f(\mathbf{x}) = W_L \sigma(W_{L-1}\sigma(W_{L-2}\sigma(\ldots W_2\sigma(W_1\mathbf{x}))))$  for a network with *L* layers. Computing  $f(\mathbf{x})$  is typically referred to as the forward pass.

We will not dwell on the details of the optimization problem in Equation (1), which is typically carried out via stochastic gradient descent algorithms or variations thereof, efficiently implemented via back-propagation (see instead, e.g., LeCun et al.,<sup>14</sup> Sun,<sup>48</sup> and Schmidhuber).<sup>49</sup> Rather, in this section we summarize some of the most popular models widely used in protein structural modeling, including how different approaches are best suited for particular data types or applications. High-level diagrams of the major architectures are shown in Figures 3.

#### **Convolutional Neural Networks**

Convolutional networks architectures<sup>50</sup> are most commonly applied to image analysis or other problems where shift-invariance or covariance is needed. Inspired by the fact that an object on an image can be shifted in the image and still be the same object, convolutional neural networks (CNNs) adopt convolutional kernels for the layer-wise affine transformation to capture this translational invariance. A 2D convolutional kernel **w** applied to a 2D image data **x** can be defined as

$$\mathbf{S}(i,j) = (\mathbf{x} * \mathbf{w})(i,j) = \sum_{m} \sum_{n} \mathbf{x}(m,n) \mathbf{w}(i-m,j-n),$$
(Equation 2)

(*i*,*i*), **x**(*m*,*n*) is the

where **S**(*i*,*j*) represents the output at position (*i*,*j*), **x**(*m*,*n*) is the value of the input **x** at position (*m*,*n*), **w**(*i*-*m*,*j*-*n*) is the parameter of kernel **w** at position (*i* - *m*,*j* - *n*), and the summation is over all possible positions. An important variant of CNN is the residual network (ResNet),<sup>51</sup> which incorporates skip-connections between layers. These modification have shown great advantages in practice, aiding the optimization of these typically huge models. CNNs, especially ResNets, have been widely used in protein structure prediction. An example is AlphaFold,<sup>22</sup> which used ResNets to predict protein inter-residue distance maps from amino acid sequences (Figure 3A).

#### **Recurrent Neural Networks**

Recurrent architectures are based on applying several iterations of the same function along a sequential input.<sup>52</sup> This can be seen as an *unfolded* architecture, and it has been widely used to process sequential data, such as time series data and written text



(i.e., NLP). With an initial hidden state  $h^{(0)}$  and sequential data  $[\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}]$ , we can obtain hidden states recursively:

Patterns Review

$$\boldsymbol{h}^{(t)} = \boldsymbol{g}^{(t)} \left( \boldsymbol{x}^{(t)}, \boldsymbol{x}^{(t-1)}, \boldsymbol{x}^{(t-2)}, \dots, \boldsymbol{x}^{(1)} \right) = f \left( \boldsymbol{h}^{(t-1)}, \boldsymbol{x}^{(t)}; \boldsymbol{\theta} \right),$$
 (Equation 3)

where *f* represents a function or transformation from one position to the next, and  $g^{(t)}$  represents the accumulative transformation up to position *t*. The hidden state vector at position *i*,  $h^{(i)}$ , contains all the information that has been seen before. As the same set of parameters (usually called a cell) can be applied recurrently along the sequential data, an input of variable length can be fed to a recurrent neural network (RNN). Due to the gradient vanishing and explosion problem (the error signal decreases or increases exponentially during training), more recent variants of standard RNNs, namely long short-term memory (LSTM)<sup>53</sup> and gated recurrent unit<sup>54</sup> are more widely used. An example of an RNN approach in the context of protein structure prediction is using an N-terminal subsequence of a protein to predict the next amino acid in the protein sequence (Figure 3B; e.g., Müller et al.<sup>55</sup>).

In conjunction with recurrent networks, attention mechanisms were first proposed (in an encoder-decoder framework) to learn which parts of a source sentence are most relevant to predicting a target word.<sup>56</sup> Compared with RNN models, attention-based models are more parallelizable and better at capturing longrange dependencies, and they are driving big advances in NLP.<sup>57,58</sup> Recently, the transformer model, which solely adopted attention layers without any recurrent or convolutional layers, was able to surpass state-of-the-art methods on language translation tasks.<sup>57</sup> For proteins, these methods could learn which parts of an amino acid sequence are critical to predicting a target residue or the properties of a target residue. For example, transformer-based models have been used to generate protein sequences conditioned on target structure,<sup>23</sup> learn protein sequence data to predict secondary structure and fitness landscapes,<sup>59</sup> and to encode the context of the binding partner in antibody-antigen binding surface prediction.<sup>60</sup>

#### Variational Autoencoder

AEs,<sup>61</sup> unlike the networks discussed so far, provide a model for unsupervised learning. Within this unsupervised

Figure 4. Different Types of Representation Schemes Applied to a Protein

framework, an AE does not learn labeled outputs but instead attempts to learn some representation of the original input. This is typically accomplished by training two parametric maps: an encoder function  $g: \mathcal{X} \to \mathbb{R}^m$ that maps an input **x** to an *m*-dimensional representation or latent space, and a decoder intended to implement the inverse map so that  $f(g(\mathbf{x})) \approx \mathbf{x}$ . Typically, the latent representation is of small dimension (*m* is smaller than

the ambient dimension of  $\mathcal{X}$ ) or constrained in some other way (e.g., through sparsity).

Variational autoencoders (VAEs),<sup>32,62</sup> in particular, provide a stochastic map between the input space and the latent space. This map is beneficial because, while the input space may have a highly complex distribution, the distribution of the representation z can be much simpler; e.g., Gaussian. These methods derive from variational inference, a method from machine learning that approximates probability densities through optimization.<sup>63</sup> The stochastic encoder, given by the *inference model*  $q_{\varphi}(\mathbf{z}|\mathbf{x})$  and parametrized by weights  $\varphi$ , is trained to approximate the true posterior distribution of the representation given the data,  $p_{\theta}(\boldsymbol{z}|\boldsymbol{x})$ . The decoder, on the other hand, provides an estimate for the data given the representation,  $p_{\theta}(\mathbf{x}|\mathbf{z})$ . Direct optimization of the resulting objective is intractable, however. Thus, training is done by maximizing the "evidence lower bound,"  $\mathcal{L}_{\theta,\varphi}(\mathbf{x})$ , instead, which provides a lower bound on the log-likehood of the data:

$$\mathcal{L}_{\theta,\varphi}(\boldsymbol{x}) = \mathsf{E}_{\boldsymbol{z} \sim q_{\varphi}(\boldsymbol{z}|\boldsymbol{x})} \ \mathsf{logp}_{\theta}(\boldsymbol{x}|\boldsymbol{z}) - \mathcal{D}_{\mathcal{KL}}(q_{\varphi}(\boldsymbol{z}|\boldsymbol{x}))||p_{\theta}(\boldsymbol{z}|\boldsymbol{x})).$$
(Equation 4)

Here,  $D_{KL}(q_{\varphi}||p_{\theta})$  is the Kullback-Leibler divergence, which quantifies the distance between distributions  $q_{\varphi}$  and  $p_{\theta}$ . Employing Gaussians for the factorized variational and likelihood distributions, as well as using a change of variables via differentiable maps, allows for the efficient optimization of these architectures.

An example of applying VAE in the protein modeling field is learning a representation of antimicrobial protein sequences (Figure 3C; e.g., Das et al.<sup>64</sup>). The resulting continuous real-valued representation can then be used to generate new sequences likely to have antimicrobial properties.

#### **Generative Adversarial Network**

Generative adversarial networks (GANs)<sup>65</sup> are another class of unsupervised (generative) models. Unlike VAEs, GANs are trained by an adversarial game between two models, or networks: a *generator*, *G*, which given a sample, *z*, from some simple distribution  $p_z(z)$  (e.g., Gaussian), seeks to map it to the distribution of some data class (e.g., naturally looking images); and a *discriminator*, *D*, whose task is to detect whether the images are real (i.e., belonging to the true distribution of the data,







Table 1. Features Cor	ntained by CUProtein Dataset			
Feature Name	Description	Dimensions	Туре	IO
AA Sequence	sequence of amino acid	n × 1	21 chars	input
PSSM	position-specific scoring matrix, a residue- wise score for motifs appearance	n × 21	real [0, 1]	input
MSA covariance	covariance matrix across homologous NA sequences	n × n	real [0, 1]	input
SS	a coarse categorized secondary structure (Q3 or Q8)	n × 1	3 or 8 chars	input
Distance matrices	pairwise distance between residues (C_{\alpha} or C_{\beta})	n × n	positive real (Å)	output
Torsion angles	variable dihedral angles for each residues $(\phi,\psi)$	n × 2	real $[-\pi, +\pi]$ (radians)	output
n, number of residues in	one protein. Data from Drori et al. <sup>78</sup>			

 $p_{data}(\mathbf{x})$ ), or fake (produced by the generator). With this gamebased setup, the generator model is trained by maximizing the error rate of the discriminator, thereby training it to "fool" the discriminator. The discriminator, on the other hand, is trained to foil such fooling. The original objective function as formulated by Goodfellow et al.<sup>65</sup> is:

$$\underset{G}{\min} \ \max_{D} V(D,G) = \mathsf{E}_{\mathbf{x} \sim \rho_{data}(\mathbf{x})}[\mathsf{log}D(\mathbf{x})] + \mathsf{E}_{\mathbf{z} \sim \rho_{\mathbf{z}}(\mathbf{z})}[\mathsf{log}(1 - D(G(\mathbf{z})))].$$
(Equation 5)

Training is performed by stochastic optimization of this differentiable loss function. While intuitive, this original GAN objective can suffer from issues, such as mode collapse and instabilities during training. The Wasserstein GAN (WGAN)<sup>66</sup> is a popular extension of GAN which introduces a Wasserstein-1 distance measure between distributions, leading to easier and more robust training.<sup>67</sup>

An example of a GAN, in the context of protein modeling is learning the distribution of protein backbone distances to generate novel protein-like folds (Figure 3D).<sup>68</sup> During training, one network *G* generates folds, and a second network *D* aims to distinguish between generated folds and fake folds.

# PROTEIN REPRESENTATION AND FUNCTION PREDICTION

One of the most fundamental challenges in protein modeling is the prediction of functionality from sequence or structure. Function prediction is typically formulated as a supervised learning problem. The property to predict can either be a protein-level property, such as a classification as an enzyme or nonenzyme,<sup>69</sup> or a residue-level property, such as the sites or motifs of phosphorylation (DeepPho)<sup>70</sup> and cleavage by proteases.<sup>71</sup> The challenging part here and in the following models is how to represent the protein. Representation refers to the encoding of a protein that serves as an input for prediction tasks or the output for generation tasks. Although a deep neural network is in principle capable of extracting complex features, a well-chosen representation can make learning more effective and efficient.<sup>72</sup> In this section, we will introduce the commonly used representations of proteins in DL models (Figure 4): sequence-based, structurebased, and one special form of representation relevant to computational modeling of proteins: coarse-grained models.

#### **Amino Acid Sequence as Representation**

As the amino acid sequence contains the information essential to reach the folded structure for most proteins,<sup>42</sup> it is widely used as an input in functional prediction and structure prediction tasks. The amino acid sequence, like other sequential data, is typically converted into one-hot encoding-based representation (each residue is represented with one high bit to identify the amino acid type and all the others low) that can be directly used in many sequence-based DL techniques.73,74 However, this representation is inherently sparse and, thus, sample-inefficient. There are many easily accessible additional features that can be concatenated with amino acid sequences providing structural, evolutionary, and biophysical information. Some widely used features include predicted secondary structure, high-level biological features, such as sub-cellular localization and unique functions,<sup>75</sup> and physical descriptors, such as AAIndex,<sup>76</sup> hydrophobicity, ability to form hydrogen bonds, charge, solvent-accessible surface area, etc. A sequence can be augmented with additional data from sequence databases, such as multiple sequence alignments (MSA) or position-specific scoring matrices (PSSMs),<sup>77</sup> or pairwise residue co-evolution features. Table 1 lists typical features as used in CU-Protein.78

#### Learned Representation from Amino Acid Sequence

Because the performance of machine learning algorithms highly depends on the features we choose, labor-intensive and domain-based feature engineering was vital for traditional machine learning projects. Now, the exceptional feature extraction ability of neural networks makes it possible to "learn" the representation, with or without giving the model any labels.<sup>72</sup> As publicly available sequence data are abundant (see Table 2), a well-learned representation that utilizes these data to capture more information is of particular interest. The class of algorithms that address the label-less learning problem fall under the umbrella of unsupervised or semi-supervised learning, which extracts information from unlabeled data to reduce the number of labeled samples needed.

P	<b>CellPress</b>
	<b>OPEN ACCESS</b>

Table 2. A Summary of Publicly Available Molecular Biology Databases							
Dataset	Description	Ν	Website				
European Bioinformatics Institute (EMBL-EBI)	a collections of wide range of datasets	-	https://www.ebi.ac.uk				
National Center for Biotechnology Information (NCBI)	a collections of biomedical and genomic databases	-	https://www.ncbi.nlm.nih.gov				
Protein Data Bank (PDB)	3D structural data of biomolecules, such as proteins and nucleic acids	~160,000	https://www.rcsb.org				
Nucleic Acid Database (NDB)	structure of nucleic acids and complex assemblies	~10,560	http://ndbserver.rutgers.edu				
Universal Protein Resource (UniProt)	protein sequence and function infromations	$\sim$ 562,000	http://www.uniprot.org/				
Sequence Read Archive (SRA)	raw sequence data from "next- generation" sequencing technologies	$\sim 3 \times 10^{16}$	NCBI database				

The most straightforward way to learn from amino acid sequence is to directly apply NLP algorithms. Word2Vec<sup>79</sup> and Doc2Vec<sup>80</sup> are groups of algorithms widely used for learning word or paragraph embeddings. These models are trained by either predicting a word from its context or predicting its context from one central word. To apply these algorithm, Asgari and Mofrad<sup>81</sup> first proposed a Word2Vec-based model called BioVec, which interprets the non-overlapping 3-mer sequence of amino acids (e.g., alanine-glutamine-lysine or AQL) as "words" and lists of shifted "words" as "sentences." They then represent a protein as the summation of all overlapping sequence fragments of length k, or k-mers (called ProtVec). Predictions based on the ProtVec representation outperformed state-of-the-art machine learning methods in the Pfam protein family<sup>82</sup> classification (93% accuracy for ~7,000 proteins, versus 69.1%-99.6%83 and 75%<sup>84</sup> for previous methods). Many Doc2Vec-type extensions were developed based on the 3-mer protocol. Yu et al.<sup>85</sup> showed that non-overlapping k-mers perform better than the overlapping ones, and Yang et al.<sup>86</sup> compared the performance of all Doc2Vec frameworks for thermostability and enantioselectivity prediction.

In these approaches, the three-residue segmentation of a protein sequence is arbitrary and does not embody any biophysical meaning. Alternatively, Alley et al.<sup>87</sup> directly used an RNN (unidirectional multiplicative long-short-term-memory or mLSTM)88 model, called UniRep, to summarize arbitrary length protein sequences into a fixed-length real representation by averaging over the representation of each residue.<sup>87</sup> Their representation achieved lower mean squared errors on 15 property prediction tasks (e.g., absorbance, activity, stability) compared with former models, including Yang et al.'s<sup>86</sup> Doc2Vec. Heinzinger et al.<sup>89</sup> adopted the bidirectional LSTM in a manner similar to Peters et al.'s<sup>90</sup> ELMo (Embeddings from Language Models) model and surpassed Asgari and Mofrad's<sup>81</sup> Word2Vec model at predicting secondary structure and regions with intrinsic disorder at the per-residue level.<sup>89</sup> The success of the transformer model in language processing, especially those trained on large number of parameters, such as BERT<sup>58</sup> and GPT3,<sup>91</sup> has inspired its application in biological sequence modeling. Rives et al.<sup>59</sup> trained a transformer model with 670 million parameters on 86 billion amino acids across 250 million protein sequences spanning evolutionary diversity. Their transformer model was superior to traditional LSTM-based models on tasks, such as the prediction of secondary structure and long-range contacts, as well as the effect of mutations on activity on deep mutational scanning benchmarks.

AEs can also provide representations for subsequent supervised tasks.<sup>32</sup> Ding et al.<sup>92</sup> showed that a VAE model is able to capture evolutionary relationships between sequences and stability of proteins, while Sinai et al.93 and Riesselman et al.94 showed that the latent vectors learned from VAEs are able to predict the effects of mutations on fitness and activity for a range of proteins, such as poly(A)-binding protein, DNA methyltransferase, and  $\beta$ -lactamase. Recently, a lower-dimensional embedding of the sequence was learned for the more complex task of structure prediction.<sup>78</sup> Alley et al.'s<sup>87</sup> UniRep surpassed former models, but since UniRep is trained on 24 million sequences and previous models (e.g., Prot2Vec) were trained on much smaller datasets (0.5 million), it is not clear if the improvement was due to better methods or the larger training dataset. Rao et al.<sup>95</sup> introduced multiple biological-relevant semi-supervised learning tasks, TAPE, and benchmarked the performance against various protein representations. Their results show conventional alignment-based inputs still outperform current selfsupervised models on multiple tasks, and the performance on a single task cannot evaluate the capacity of models. A comprehensive and persuasive comparison of representations is required.

#### **Structure as Representation**

Since the most important functions of a protein (e.g., binding, signaling, catalysis) can be traced back to the 3D structure of the protein, direct use of 3D structural information, and analogously, learning a good representation based on 3D structure, are highly desired. The direct use of raw 3D representations (such as coordinates of atoms) is hindered by considerable challenges, including the processing of unnecessary information due to translation, rotation, and permutation of atomic indexing. Townshend et al.<sup>96,97</sup> and Simonovsky and Meyers<sup>96,97</sup> obtained a translationally invariant, 3D representation of each residue by voxelizing its atomic neighborhood for a grid-based 3D CNN model. The work of Kolodny et al.,<sup>98</sup> Taylor,<sup>99</sup> and Li and



Koehl<sup>100</sup> representing the 3D structure of a protein as 1D strings of geometric fragments for structure comparison and fold recognition may also prove useful in DL approaches. Alternatively, the torsion angles of the protein backbone, which are invariant to translation and rotation, can fully recapitulate protein backbone structure under the common assumption that variation in bond lengths and angles is negligible. AlQuraishi<sup>101</sup> used backbone torsion angles to represent the 3D structure of the protein as a 1D data vector. However, because a change in a backbone torsion angle at a residue affects the inter-residue distances between all preceding and subsequent residues, these 1D variables are highly interdependent, which can frustrate learning. To circumvent these limitations, many approaches use 2D projections of 3D protein structure data, such as residue-residue distance and contact maps,<sup>24,102</sup> and pseudo-torsion angles and bond angles that capture the relative orientations between pairs of residues.<sup>103</sup> While these representations guarantee translational and rotational invariance, they do not guarantee invertibility back to the 3D structure. The structure must be reconstructed by applying constraints on distance or contact parameters using algorithms, such as gradient descent minimization, multidimensional scaling, a program like the Crystallography and NMR system (CNS),<sup>104</sup> or in conjunction with an energy-function-based protein structure prediction program.<sup>22</sup>

An alternative to the above approaches for representing protein structures is the use of a graph, i.e., a collection of nodes or vertices connected by edges. Such a representation is highly amenable to the graph neural network (GNN) paradigm, <sup>105</sup> which has recently emerged as a powerful framework for non-Euclidean data<sup>106</sup> in which the data are represented with relationships and inter-dependencies, or edges between objects or nodes.<sup>107</sup> While the representation of proteins as graphs and the application of graph theory to study their structure and properties has a long history,<sup>108</sup> the efforts to apply GNNs to protein modeling and design is quite recent. As a benchmark, many GNNs<sup>69,109</sup> have been applied to classify enzymes from non-enzymes in the PROTEINS<sup>110</sup> and D&D<sup>111</sup> datasets. Fout et al.<sup>112</sup> utilized a GNN in developing a model for protein-protein interface prediction. In their model, the node feature comprised residue composition and conservation, accessible surface area, residue depth, and protrusion index; and the edge feature comprised a distance and an angle between the normal vectors of the amide plane of each node/residue. A similar framework was used to predict antibody-antigen binding interfaces.<sup>60</sup> Zamora-Resendiz and Crivelli<sup>113</sup> and Gligorijevic et al.<sup>114</sup> further generalized and validated the use of graph-based representations and the graph convolutional network (GCN) framework in protein function prediction tasks, using a class activation map to interpret the structural determinants of the functionalities. Torng and Altman<sup>115</sup> applied GCNs to model pocket-like cavities in proteins to predict the interaction of proteins with small molecules, and Ingraham et al.<sup>23</sup> adopted a graph-based transformer model to perform a protein sequence design task. These examples demonstrate the generality and potential of the graph-based representation and GNNs to encode structural information for protein modeling.

The surface of the protein or a cavity is an information-rich region that encodes how a protein may interact with other molecules and its environment. Recently, Gainza et al.<sup>116</sup> used a geometric DL framework<sup>117</sup> to learn a surface-based representation of the protein, called MaSIF. They calculated "fingerprints" for patches on the protein surface using geodesic convolutional layers, which were further used to perform tasks, such as binding site prediction or ultra-fast protein-protein interaction (PPI) search. The performance of MaSIF approached the baseline of current methods in docking and function prediction, providing a proof-of-concept to inspire more applications of geometrybased representation learning.

#### **Score Function and Force Field**

A high-quality force field (or, more generally, score function) for sampling and/or ranking models (decoys) is one of the most vital requirements for protein structural modeling.<sup>118</sup> A force field describes the potential energy surface of a protein. A score function may contain knowledge-based terms that do not necessarily have a valid physical meaning, and they are designed to distinguish near-native conformations from non-native ones (for example, learning the GDT\_TS).<sup>119</sup> A molecular dynamics (MD) or Monte Carlo (MC) simulation with a state-of-the-art force field or score function can reproduce reasonable statistical behaviors of biomolecules.<sup>120–122</sup>

Current DL-based efforts to learn the force field can be divided into two classes: "fingerprint" based and graph based. Behler and Parrinello<sup>123</sup> developed roto-translationally invariant features, i.e., the Behler-Parrinello fingerprint, to encode the atomic environment for neural networks to learn potential surfaces from density functional theory (DFT) calculations. Smith et al. extended this framework and tested its accuracy by simulating systems up to 312 atoms (Trp-cage) for 1 ns.<sup>124,125</sup> Another family that includes deep tensor neural networks<sup>126</sup> and SchNet<sup>127</sup> uses graph convolutions to learn a representation for each atom within its chemical environment. Although the prediction quality and the ability to learn a representation with novel chemical insight make the graph-based approach increasingly popular,<sup>28</sup> the application scales poorly to larger systems and thus has mainly focused on small organic molecules.

We anticipate a shift toward DL-based score functions because of the enormous gains in speed and efficiency. For example, Zhang et al.<sup>128</sup> showed that MD simulation on a neural potential was able to reproduce energies, forces, and time-averaged properties comparable with ab initio MD (AIMD) at a cost that scales linearly with system size, compared with cubic scaling typical for AIMD with DFT. Although these force fields are, in principle, generalizable to larger systems, direct applications of learned potentials to model full proteins are still rare. PhysNet, trained on a set of small peptide fragments (at most eight heavy atoms), was able to generalize to deca-alanine (Ala10),<sup>129</sup> and ANI-1x and AIMNet have been tested on chignolin (10 residues) and Trp-cage (20 residues) within the ANI-MD benchmark dataset.<sup>125,130</sup> Lahey and Rowley<sup>131</sup> and Wang et al.<sup>132</sup> combined the quantum mechanics/molecular mechanics (QM/MM) strategy<sup>133</sup> and the neural potential to model docking with small ligands and larger proteins.<sup>131,132</sup> Recently, Wang et al.<sup>134</sup> proposed an end-to-end differential MM force field by training a GNN on energies and forces to learn atom-typing and force field parameters.

#### **Coarse-Grained Models**

Coarse-grained models are higher-level abstractions of biomolecules, such as using a single pseudo-atom or a bead to

represent multiple atoms, grouped based on local connectivity and/or chemical properties. Coarse graining smoothens out the energy landscape, and thereby helps avoid trapping in local minima and speeds up conformational sampling.<sup>135</sup> One can learn the atomic-level properties to construct a fast and accurate neural coarse-grained model once the coarse-grained mapping is given. Early attempts to apply DL-based methods to coarsegraining focus on water molecules with the roto-translationally invariant features.<sup>136,137</sup> Wang et al.<sup>138</sup> developed CGNet and learned the coarse-grained model of the mini protein, chignolin, in which the atoms of a residue are mapped to the corresponding  $C_{\alpha}$  atom. The free energy surface learned with CGNet is quantitatively correct and MD simulations performed with CGNet potentially predict the same set of metastable states (folded, unfolded, and misfolded). Another critical question for coarse graining is determining which sets of atoms to map into a united atom. For example, one choice is to use a single coarse-grained atom to represent a whole residue, and a different choice is to use two coarse-grained atoms, one to represent the backbone and the other to represent the side chain. To determine the optimal choice, Wang and Gómez-Bombarelli<sup>139</sup> applied an encoder-decoder-based model to explicitly learn the lowerdimensional representation of proteins by minimizing the information loss at different levels of coarse graining. Li et al.140 treated this problem as a graph segmentation problem and presented a GNN-based coarse-graining mapping predictor called Deep Supervised Graph Partitioning Model.

#### **STRUCTURE DETERMINATION**

The most successful application of DL in the field of protein modeling so far has been the prediction of protein structure. Protein structure prediction is formulated as a well-defined problem with clear inputs and outputs: predict the 3D structure (output) given amino acid sequences (input), with the experimental structures as the ground truth (labels). This problem perfectly fits the classical supervised learning approach, and once the problem is defined in these terms, the remaining challenge is to choose a framework to handle the complex relationship between input and output. The CASP experiment for structure prediction is held every 2 years and served as a platform for DL to compete with state-of-the-art methods and, impressively, outshine them in certain categories. We will first discuss the application of DL to the protein folding problem, and then comment on some problems related to structure determination. Table 3 summarizes major DL efforts in structure prediction.

#### **Protein Structure Prediction**

Before the notable success of DL at CASP12 (2016) and CASP13 (2018), the state-of-the-art methodology used complex workflows based on a combination of fragment insertion and structure optimization methods, such as simulated annealing with a score function or energy potential. Over the last decade, the introduction of co-evolution information in the form of evolutionary coupling analysis (ECA)<sup>154</sup> improved predictions. ECA relies on the rationale that residue pairs in contact in 3D space tend to evolve or mutate together; otherwise, they would disrupt the structure to destabilize the fold or render a large conformational change. Thus, evolutionary couplings from sequencing data suggest distance relationships between residue pairs and aid structure construction from sequence through contact or distance constraints. Because co-evolution information relies on statistical averaging of sequence information from a large number of MSAs, 145, 155, 156 this approach is not effective when the protein target has only a few sequence homologs. Neural networks were, at first, introduced to deduce evolutionary couplings between distant homologs, thereby improving ECA-type contact predictions for contact-assisted protein folding.154 While the application of neural networks to learn inter-residue protein contacts dates back to the early 2000s,<sup>157,158</sup> more recently this approach was adopted by MetaPSICOV (two-layer NN),146 PConsC2 (two-layer NN),145 and CoinDCA-NN (five-layer NN),<sup>155</sup> which combined neural networks with ECAs. However, there was no significant advantage to neural networks compared with other machine learning methods at that time.<sup>159</sup>

In 2017, Wang et al.<sup>102</sup> proposed RaptorX-Contact, a residual neural network (ResNet)-based model,<sup>51</sup> which, for the first time used a deep neural network for protein contact prediction, significantly improving the accuracy on blind, challenging targets with novel folds. RaptorX-Contact ranked first in free modeling targets at CASP12.<sup>161</sup> Its architecture (Figure 5(a)) entails (1) a 1D ResNet that inputs MSAs, predicted secondary structure and solvent accessibility (from DL-based prediction tool RaptorX-Property)<sup>162</sup> and (2) a 2D ResNet with dilations that inputs the 1D ResNet output and inter-residue co-evolution information from CCMpred.<sup>144</sup> In its original formulation, RaptorX-Contact outputs a binary classification of contacting versus non-contacting residue pairs.<sup>102</sup> Later versions were trained to learn multiclass classification for distance distributions between C<sub>B</sub> atoms.<sup>147</sup> The primary contributors to the accuracy of predictions was the co-evolution information from CCMpred and the depth of the 2D ResNet, suggesting that the deep neural network learned co-evolution information better than previous methods. Later, the method was extended to predict  $C_{\alpha}-C_{\alpha}$ ,  $C_{\alpha}-C_{\gamma}$ ,  $C_{\gamma}$ - $C_{\gamma}$ , N-O distances and torsion angles (DL-based RaptorX-Angle),<sup>163</sup> giving constraints to locate side chains and additionally constrain the backbone; all five distances, torsions, and secondary structure predictions were converted to constraints for folding by CNS.<sup>147</sup> At CASP12, however, RaptorX-Contact (original contact-based formulation) and DL drew limited attention because the difference between top-ranked predictions from DL-based methods and hybrid DCA-based methods was small.

This situation changed at CASP13 4 when one DL-based model, AlphaFold, developed by team A7D, or Deep-Mind,<sup>26,22,164</sup> ranked first and significantly improved the accuracy of "free modeling" (no templates available) targets (Figure 1). The A7D team modified the traditional simulated annealing protocol with DL-based predictions and tested three protocols based on deep neural networks. Two protocols used memory-augmented simulated annealing (with domain segmentation and fragment assembly) with potentials generated from predicted inter-residue distance distributions and predicted GDT\_TS,<sup>165</sup> respectively, whereas the third protocol directly applies gradient descent optimization on a hybrid potential combining predicted distance and Rosetta score. For the distance prediction network, a deep ResNet, similar to that of RaptorX,<sup>102</sup> inputs MSA data and predicts the probability of distances between  $\beta$ - carbons. A second network was

## CellPress OPEN ACCESS



Table 3. A Summary of Structure Prediction Models							
Model	Architecture	Dataset	N_train	Performance	Testset	Citation	
/	MLP(2-layer)	proteases	13	3.0 Å RMSD (1TRM),1.2 Å RMSD (6PTI)	1TRM, 6PTI	Bohr et al. <sup>9</sup>	
PSICOV	graphical Lasso	-	-	precision: Top-L 0.4, Top-L/2 0.53,Top-L/5 0.67, Top-L/ 10 0.73	150 Pfam	Jones et al. <sup>141</sup>	
CMAPpro	2D biRNN + MLP	ASTRAL	2,352	precision: Top-L/5 0.31, Top-L/ 10 0.4	ASTRAL 1.75 CASP8, 9	Di Lena et al. <sup>142</sup>	
DNCON	RBM	PDB SVMcon	1,230	precision: Top-L 0.46, Top-L/2 0.55, Top-L/5 0.65	SVMCON_TEST, D329, CASP9	Eickholt et al. <sup>143</sup>	
CCMpred	LM	-	-	precision: Top-L 0.5, Top-L/2 0.6, Top-L/5 0.75, Top-L/10 0.8	150 Pfam	Seemayer et al. <sup>144</sup>	
PconsC2	Stacked RF	PSICOV set	150	positive predictive value (PPV) 0.44	set of 383 CASP10(114)	Skwark et al. <sup>145</sup>	
MetaPSICOV	MLP	PDB	624	precision: Top-L 0.54, Top-L/2 0.70, Top-L/5 0.83, Top-L/ 10 0.88	150 Pfam	Jones et al. <sup>146</sup>	
RaptorX-Contact	ResNet	subset of PDB25	6,767	TM score: 0.518 (CCMpred: 0.333, MetaPSICOV: 0.377)	Pfam, CASP11, CAMEO, MP	Wang et al, 2017 <sup>102</sup>	
RaptorX-Distance	ResNet	subset of PDB25	6,767	TM score: 0.466 (CASP12), 0.551 (CAMEO), 0.474 (CASP13)	CASP12 + 13, CAMEO	Xu, 2018 <sup>147</sup>	
DeepCov	2D CNN	PDB	6,729	precision: Top-L 0.406, Top-L/2 0.523, Top-L/5 0.611, Top-L/ 10 0.642	CASP12	Jones et al, 2018 <sup>148</sup>	
SPOT	ResNet, Res-bi-LSTM	PDB	11,200	AUC: 0.958 (RaptorX-contact ranked 2nd: 0.909)	1,250 chains after June 2015	Hanson et al. <sup>149</sup>	
DeepMetaPSICOV	ResNet	PDB	6,729	precision: Top-L/5 0.6618	CASP13	Kandathil et al, 2019 <sup>150</sup>	
MULTICOM	2D CNN	CASP 8-11	425	TM score: 0.69, GDT_TS: 63.54, SUM <i>Z</i> score (– 2.0): 99.47	CASP13	Hou et al. <sup>151</sup>	
C-I-TASSER*	2D CNN	-	-	TM score: 0.67, GDT_HA: 0.44, RMSD: 6.19, SUM <i>Z</i> score( – 2.0): 107.59	CASP13	Zheng et al. <sup>152</sup>	
AlphaFold	ResNet	PDB	31,247	TM score: 0.70, GDT_TS: 61.4,SUM <i>Z</i> score (– 2.0): 120.43	CASP13	Senior et al. <sup>22</sup>	
MapPred	ResNet	PISCES	7,277	precision: 78.94% in SPOT, 77.06% in CAMEO, 77.05 in CASP12	SPOT, CAMEO, CASP12	Wu et al, 2019 <sup>153</sup>	
trRosetta	ResNet	PDB	15,051	TM_score: 0.625 (AlphaFold: 0.587)	CASP13, CAMEO	Yang et al, 2020 <sup>103</sup>	
RGN	bi-LSTM	ProteinNet 12 (before 2016)**	104,059	10.7 Å dRMSD on FM, 6.9 Å on TBM	CASP12	AlQuraishi, 2019 <sup>101</sup>	
/	biGRU, Res LSTM	CUProtein	75,000	preceded CASP12 winning team, comparable with AlphaFold in RMSD	CASP12 + 13	Drori et al. <sup>78</sup>	

FM, free modeling; GRU, gated recurrent unit; LM, pseudo-likelihood maximization; MLP, multi-layer perceptron; MP, membrane protein; RBM, restricted Boltzmann machine; RF, random forest; RMSD, root-mean square deviation; TBM, template-based modeling.

\*C-I-TASSER and C-QUARK were reported, we only report one here.

\*\*RGN was trained on different ProteinNet for each CASP, we report the latest one here.

trained to predict GDT\_TS of the candidate structure with respect to the true or native structure. The simulated annealing process was improved with a conditional variational autoencoder (CVAE)<sup>166</sup> model that constructs a mapping between the backbone torsions and a latent space conditioned by sequence. With this network, the team generated a database

of nine-residue fragments for the memory-augmented simulated annealing system. Gradient-based optimization performed slightly better than the simulated annealing, suggesting that traditional simulated annealing is no longer necessary and state-of-the-art performance can be reached with simply optimizing a network predicted potential. AlphaFold's authors,





#### Figure 5. Two Representative DL Approaches to Protein Structure Prediction

(A) Residue distance prediction by RaptorX: the overall network architecture of the deep dilated ResNet used in CASP13. Inputs of the first-stage, 1D convolutional layers are a sequence profile, predicted secondary structure, and solvent accessibility. The output of the first stage is then converted into a 2D matrix by concatenation and fed into a deep ResNet along with pairwise features (co-evolution information, pairwise contact, and distance potential). A discretized interresidue distance is the output. Additional network layers can be attached to predict torsion angles and secondary structures. Figure from Xu and Wang (2019).<sup>160</sup>
(B) Direct structure prediction: overview of recurrent geometric networks (RGN) approach. The raw amino acid sequence along with a PSSM are fed as input features, one residue at a time, to a bidirectional LSTM net. Three torsion angles for each residue are predicted to directly construct the 3D structure. Figure from Xu and AlQuraishi (2019).<sup>101</sup>

like the RaptorX-Contact group, emphasized that the accuracy of predictions relied heavily on learned distance distributions and co-evolutionary data.

Yang et al.<sup>103</sup> further improved the accuracy of predictions on CASP13 targets using a shallower network than former models (61 versus 220 ResNet blocks in AlphaFold) by also training their neural network model (named trRosetta) to learn inter-residue

orientations along with  $\beta-$  carbon distances. The geometric features— $C_{\alpha}$ - $C_{\beta}$  torsions, pseudo-bond angles, and azimuthal rotations—directly describe the relevant coordinates for the physical interaction of two amino acid side chains. These additional outputs created significant improvement on a relatively fixed DL framework, suggesting that there is room for additional improvement.



An alternative and intuitive approach to structure prediction is directly learning the mapping from sequence to structure with a neural network. AlQuraishi<sup>101</sup> developed such an end-to-end differentiable protein structure predictor, called RGN, that allows direct prediction of torsion angles to construct the protein backbone (Figure 5B). RGN is a bidirectional LSTM that inputs a sequence, PSSM, and positional information and outputs predicted backbone torsions. Overall 3D structure predictions are within 1–2 Å of those made by top-ranked groups at CASP13, and this approach boasts a considerable advantage in prediction time compared with strategies that learn potentials. Moreover, the method does not use MSA-based information and could potentially be improved with the inclusion of evolutionary information. The RGN strategy is generalizable and well suited for protein structure prediction. Several generative methods (see below) also entail end-to-end structure prediction models, such as the CVAE framework used by AlphaFold, albeit with more limited success.<sup>22</sup>

#### **Related Applications**

Side-chain prediction is required for homology modeling and various protein engineering tasks, such as fixed-backbone design. Side-chain prediction is often embedded in high-resolution structure prediction methods, traditionally with dead-end elimination<sup>167</sup> or preferential sampling from backbone-dependent side-chain rotamer libraries.<sup>168</sup> Liu et al.<sup>169</sup> specifically trained a 3D CNN to evaluate the probability score for different potential rotamers. Du et al.<sup>170</sup> adopted an energy-based model (EBM)<sup>171</sup> to recover rotamers for backbone structures. Recent protein structure prediction models, such as Gao et al.'s<sup>163</sup> RaptorX-angle and Yang et al.'s<sup>103</sup> trRosetta, predict the structural features that help locate the position of side-chain atoms as well.

PPI prediction identifies residues at the interface of the two proteins forming a complex. Once the interface residues are determined, a local search and scoring protocol can be used to determine the structure of a complex. Similar to protein folding, efforts have focused on learning to classify contact or not. For example, Townshend et al.96 developed a 3D CNN model (SASNet) that voxelizes the 3D environment around the target residue, and Fout et al.<sup>112</sup> developed a GCN-based model with each interacting partner represented as a graph. Unlike those starting from the unbound structures, Zeng et al.<sup>172</sup> reuse the model trained on single-chain proteins (i.e., RaptorX-Contact) to predict PPI with sequence information alone, which resulted in the RaptorX-Complex that outperforms ECA-based methods at contact prediction. Another interesting approach directly compares the geometry of two protein patches. Gainza et al.<sup>116</sup> trained their MaSIF model by minimizing the Euclidean distances between the complementary surface patches on the two proteins while maximizing the distances between non-interacting surface patches. This step is followed by a quick nearestneighbor scanning to predict binding partners. The accuracy of MaSIF was comparable with traditional docking methods. However, MaSIF, similar to existing methods, showed low prediction accuracy for targets that involve conformational changes during binding.

Membrane proteins (MPs) are partially or fully embedded in a hydrophobic environment composed of a lipid bilayer and, consequently, they exhibit hydrophobic motifs on the surface



unlike the majority of the proteins that are water soluble. Wang et al.<sup>173</sup> used a DL transfer learning framework comprising one-shot learning from non-MPs to MPs. They showed that transfer learning works surprisingly well here because the most frequently occurring contact patterns in soluble proteins and MPs are similar. Other efforts include classification of the trans-membrane topology.<sup>174</sup> Since experimental biophysical data are sparse for MPs, Alford and Gray<sup>175</sup> compiled a collection of 12 diverse benchmark sets for membrane protein prediction and design for testing and learning of implicit membrane energy models.

Loop modeling is a special case of structure prediction, where most of the 3D protein structure is given, but coordinates of segments of the polypeptide are missing and need to be completed. Loops are irregular and sometimes flexible segments, and thus their structures have been difficult to capture experimentally or computationally.<sup>176,177</sup> So far, DL frameworks based on interresidue distance prediction (similar to protein structure prediction),<sup>178</sup> and those based on treating the loop residue distances with the remaining residues as an image inpainting problem<sup>179</sup> have been applied to loop modeling. Recently, Ruffolo et al.<sup>177</sup> used a RaptorX-like network setup and a trRosetta geometric representation to predict the structure of antibody hypervariable complementarity-determining region (CDR) H3 loops, which is critical for antigen binding.

#### **PROTEIN DESIGN**

We divide the current DL approaches to protein design into two broad categories. The first uses knowledge of other sequences (either "all" sequenced proteins or a certain class of proteins) to design sequences directly (Table 4). These approaches are well suited to create new proteins with functionality matching existing proteins based on sequence information alone, in a manner similar to consensus design.<sup>180</sup> The second class follows the "fold-before-function" scheme and seeks to stabilize specific 3D structures, perhaps but not necessarily with the intent to perform a desired function (Tables 5 and 6). The first approach can be described as function  $\rightarrow$  sequence (structure agnostic), and the second approach fits the traditional stepwise inverse design: function  $\rightarrow$  structure  $\rightarrow$  sequence.

Many of the recent studies describe novel algorithms that output putative designed protein sequences, but only a few studies also present experimental validation. In traditional protein design studies, it is not uncommon for most designs to fail, and some of the early reports of protein designs were later withdrawn when the experimental evidence was not confirmed by others. As a result, it is usually expected that design studies offer rigorous experimental evidence. In this review, because we are interested in creative, emerging DL methods for design, we include papers that lack experimental validation, and many of these have *in silico* tests that help gauge validity. In addition, we make a special note of recent studies that present experimental validation of designs.

#### **Direct Design of Sequence**

Approaches that attempt to design for sequences parallel work in the field of NLP, where an auto-regressive framework is common, most notably, the RNN. In language processing, an RNN



Table 4. Generative Models to Identify Sequence from Function (Design for Function)						
Model	Architecture	Output	Dataset	N_train	Performance	Citation
-	WGAN + AM	DNA	chromosome 1 of human hg 38	4.6M	~4 times stronger than training data in predicted TF binding	Killoran et al. <sup>181</sup>
-	VAE	AA	5 protein families	-	natural mutation probability prediction rho = 0.58	Sinai et al. <sup>93</sup>
-	LSTM	AA	ADAM, APD, DADP	1,554	predicted antimicrobial property $0.79 \pm 0.25$ (random: $0.63 \pm 0.26$ )	Müller et al, 2018 <sup>55</sup>
PepCVAE	CVAE	AA	-	15K labeled, 1.7M unlabeled	generate predicted AMP with 83% (random, 28%; length, 30)	Das et al. <sup>64</sup>
FBGAN	WGAN	DNA	UniProt (res., 50)	3,655	predicted antimicrobial property over 0.9 after 60 epochs	Gupta et al. <sup>182</sup>
DeepSequence	VAE	AA	mutational scan data	41 scans	aimed for mutation effect prediction, outperformed previous models	Riesselman et al. <sup>94</sup>
DbAS-VAE	VAE+AS	DNA	simulated data	-	predicted protein expression surpassed FB-GAN/VAE	Brookes et al. <sup>183</sup>
-	LSTM	musical scores	-	56 betas + 38 alphas	generated proteins capture the secondary structure feature	Yu et al <sup>.184</sup>
BioSeqVAE	VAE	AA	UniProt	200,000	83.7% reconstruction accuracy,70.6% EC accuracy	Costello et al. <sup>185</sup>
-	WGAN	AA	antibiotic resistance determinants	6,023	29% similar to training sequence (BLASTp)	Chhibbar et al. <sup>186</sup>
PEVAE	VAE	AA	3 protein families	31,062	latent space captures phylogenetic, ancestral relationship, and stability	Ding et al. <sup>92</sup>
-	ResNet	AA	mutation data + Ilama immune repertoire	1.2M (nano)	predicted mutation effect reached state-of-the-art, built a library of CDR3 seq	Riesselman et al. <sup>187</sup>
Vampire	VAE	AA	immuneACCESS	-	generated sequences predicted to be similar to real CDR3 sequences	Davidson et al, 2019 <sup>188</sup>
ProGAN	CGAN	AA	eSol	2,833	solubility prediction <i>R</i> <sup>2</sup> improved from 0.41 to 0.45	Han et al, 2019 <sup>189</sup>
ProteinGAN	GAN	AA	MDH from UniProt	16,706	60 sequences were tested <i>in vitro</i> , 19 soluble, 13 with catalytic activity	Repecka et al. <sup>190</sup>
CbAS-VAE	VAE+AS	AA	protein fluorescence dataset	5,000	predicted protein fluorescence surpassed FB-VAE/DbAS	Brookes et al. <sup>183</sup>

AA, amino acid sequence; AM, activation maximization; AS, adaptive sampling; CGAN, conditional generative adversarial network; CVAE, conditional variational autoencoder; DNA, DNA sequence; EC, enzyme commission.

model is able to take the beginning of a sentence and predict the next word in that sentence. Likewise, given a starting amino acid residue or a sequence of residues, a protein design model can output a categorical distribution for each of the 20 amino acid residues for the next position in the sequence. The next residue in the sequence is sampled from this categorical distribution, which in turn is used as the input to predict the following one. Following this approach, new sequences, sampled from the distribution of the training data, are generated, with the goal of having properties similar to those in the training set. Müller et al.<sup>55</sup> first applied an LSTM RNN framework to learn sequence patterns of antimicrobial peptides (AMPs),<sup>204</sup> a highly specialized sequence space of cationic, amphipathic helices. The same group then applied this framework to design membranolytic anticancer peptides.<sup>205</sup> Twelve of the generated peptides were synthesized and six of them killed MCF7 human breast adenocarcinoma cells with at least 3-fold selectivity against human erythrocytes. In another application, instead of traditional

RNNs, Riesselman et al.<sup>187</sup> used a residual causal dilated CNN<sup>206</sup> in an auto-regressive way and generated a functional single-domain antibody library conditioned on the naive immune repertoires from llamas; although experimental validation was not presented. Such applications could potentially speed up and simplify the task of generating sequence libraries in the lab.

Another approach to sequence generation is mapping the latent space to the sequence space, and common strategies to train such a mapping include AEs and GANs. As mentioned earlier, AEs are trained to learn a bidirectional mapping between a discrete design space (sequence) and a continuous real-valued space (latent space). Thus, many applications of AEs use the learnt latent representation to capture the sequence distribution of a specific class of proteins, and subsequently, to predict the effect of variations in sequence (or mutations) on protein function.<sup>92–94</sup> The utility of this learned latent space, however, is more than that. A well trained real-valued latent space can be used to interpolate between two training samples, or even





Table 5. Generative Models for Protein Structure Design						
Model	Architecture	Representation	Dataset	N_train	Performance	Citation
-	DCGAN	$C_{\alpha}$ - $C_{\alpha}$ distances	PDB (16-, 64-, 128-residue fragments)	115,850	meaningful secondary structure, reasonable Ramachandran plot	Anand et al. <sup>24</sup>
RamaNet	GAN	torsion angles	ideal helical structures from PDB	607	generated torsions are concentrated around helical region	Sabban et al. <sup>191</sup>
-	DCGAN	backbone distance	PDB (64-residue fragment)	800,000	smooth interpolations; recover from sequence design and folding	Anand et al. <sup>68</sup>
lg-VAE	VAE	coordinates and backbone distance	AbDb (antibody structure)	10,768	sampled 5,000 Igs screened for SARS-CoV2 Binder	Eguchi et al. <sup>192</sup>
-	CNN (input design)	same as trRosetta	-	-	27 out of 129 sequence-structure pairs experimentally validated	Anishchenko et al. <sup>193</sup>

CNN, convolutional neural network; DCGAN, deep convolutional generative adversarial network; GAN, generative adversarial network; VAE, variational autoencoder.

extrapolate beyond the training data to yield novel sequences. One such example is the PepCVAE model.<sup>64</sup> Following a semisupervised learning approach, Das et al.<sup>64</sup> trained a VAE model on an unlabeled dataset of  $1.7 \times 10^6$  sequences and then refined the model for the AMP subspace using a 15,000 sequencelabeled dataset. By concatenating a conditional code indicating if a peptide is antimicrobial, the CVAE framework allows efficient sampling of AMPs selectively from the broader peptide space. More than 82% of the generated peptides were predicted to exhibit antimicrobial properties according to a state-of-the-art AMP classifier.

Unlike AEs, GANs focus on learning the unidirectional mapping from a continuous real-valued space to the design space. In an early example, Killoran et al.'s<sup>181</sup> developed a model that combines a standard GAN and activation maximization to design DNA sequences that bind to a specific protein. Repecka et al.<sup>190</sup> trained ProteinGAN on the bacterial enzyme malate dehydrogenase (MDH) to generate new enzyme sequences that were active and soluble in vitro, some with over 100 mutations, with a 24% success rate. Another interesting GAN-based framework is Gupta and Zou's<sup>207</sup> FeedBack GAN (FBGAN) that learns to generate cDNA sequences for peptides. They add a feedbackloop architecture to optimize the synthetic gene sequences for desired properties using an oracle (an external function analyzer). At every epoch, they update the positive training data for the discriminator with high-scoring sequences from the generator so that the score of generated sequences increases gradually. They demonstrated the efficacy of their model by successfully biasing generated sequences toward antimicrobial activity and a desired secondary structure.

#### **Design with Structure as Intermediate**

Within the fold-before-function scheme, for design one first picks a protein fold or topology according to certain desirable properties, then determines the amino acid sequence that could fold into that structure (function  $\rightarrow$  structure  $\rightarrow$  sequence). Under the supervised learning setting, most efforts use the native sequences as the ground truth and recovery rate of native sequences (i.e., the percentage of sequence that matches the

native one) as a success metric. To compare, Kuhlman and Baker<sup>208</sup> reported sequence recovery rates of 51% for core residues and 27% for all amino acid residues using traditional *de novo* design approaches. Because the mapping from sequence to structure is not unique (within a neighborhood of each structure), it is not clear that higher sequence recovery rates would be meaningful.

A class of efforts, pioneered by the SPIN model,<sup>209</sup> inputs a five-residue sliding window to predict the amino acid probabilities for the center position to generate sequences compatible with a desired structure. The features in such models include o and  $\psi$  dihedrals, a sequence profile of a five-residue fragment derived from similar structures, and a rotamer-based energy profile of the target residue using the DFIRE potential. SPIN<sup>209</sup> reached a 30.7% sequence recovery rate and Wang et al.<sup>194</sup> and O'Connell et al.'s<sup>25</sup> SPIN2 further improved it to 34%. Another class of efforts inputs the voxelized local environment of an amino acid residue. In Zhang et al.'s<sup>197,198</sup> and Shroff et al.'s<sup>197,198</sup> models, voxelized local environment was fed into a 3D CNN framework to predict the most stable residue type at the center of a region. Shroff et al.<sup>198</sup> reported a 70% recovery rate and the mutation sites were validated experimentally. Anand et al.<sup>202</sup> trained a similar model to design sequences for a given backbone. Their protocol involves iteratively sampling from predicted conditional distributions, and it recovered from 33% to 87% of native sequence identities. They tested their model by designing sequences for five proteins, including a de novo TIM barrel. The designed sequences were 30%-40% identical to native sequences and predicted structures were 2-5 Å root-mean-square deviation from the native conformation.

Other approaches generate full sequences conditioned by a target structure. Greener et al.<sup>195</sup> trained a CVAE model to generate sequences conditioned on protein topology represented in a string.<sup>99</sup> The resulting sequence was verified to be stable with molecular simulation. Karimi et al.<sup>210</sup> developed gcWGAN that combined a CGAN and a guidance strategy to bias the generated sequences toward a desired structure. They used a fast structure prediction algorithm<sup>211</sup> as an "oracle"



Table 6. Generative Models to Identify Sequence from Structure (Protein Design)							
Model	Architecture	Input	Dataset	N_train	Performance	Citation	
SPIN	MLP	sliding window with 136 features	PISCES	1,532	sequence recovery of 30.7% on 1,532 proteins (CV)	Li et al. <sup>100</sup>	
SPIN2	MLP	sliding window with 190 features	PISCES	1,532	sequence recovery of 34.4% on 1,532 proteins (CV)	O'Connell et al. <sup>25</sup>	
_	MLP	target residue and its neighbor as pairs	PDB	10,173	sequence recovery of 34% on 10,173 proteins	Wang et al. <sup>194</sup>	
_	CVAE	string encoded structure or metal	PDB, MetalPDB	3,785	verified with structure prediction and dynamic simulation	Greener et al. <sup>195</sup>	
SPROF	Bi-LSTM + 2D ResNet	112 1-D features + $C_{\alpha}$ distance map	PDB	11,200	sequence recovery of 39.8% on protein	Chen et al. <sup>196</sup>	
ProDCoNN	3D CNN	gridded atomic coordinates	PDB	17,044	sequence recovery of 42.2% on 5,041 proteins	Zhang et al. <sup>197</sup>	
-	3D CNN	gridded atomic coordinates	PDB-REDO	19,436	sequence recovery 70%, experimental validation of mutation	Shroff et al. <sup>198</sup>	
ProteinSolver	Graph NN	partial sequence, adjacency matrix	UniParc	72×10 <sup>6</sup> residues	sequence recovery of 35%, folding and MD test with 4 proteins	Strokach et al, 2019 <sup>199</sup>	
gcWGAN	CGAN	random noise + structure	SCOPe	20,125	diversity and TM score of prediction from designed sequence $\geq$ cVAE	Karimi et al. <sup>200</sup>	
-	Graph Transformer	backbone structure in graph	CATH based	18,025	perplexity: 6.56 (rigid), 11.13 (flexible) (random: 20.00)	Ingraham et al. <sup>23</sup>	
DenseCPD	ResNet	gridded backbone atomic density	PISCES	2.6×10 <sup>6</sup> residues	sequence recovery of 54.45% on 500 proteins	Qi et al. <sup>201</sup>	
-	3D CNN	gridded atomic coordinates	PDB	21,147	sequence recovery from 33% to 87%, test with folding of TIM barrel	Anand et al. <sup>202</sup>	
-	CNN (input design)	Same as trRosetta	-	-		Norn et al. <sup>203</sup>	

Bi-LSTM, bidirectional long short-term memory; CV, cross-validation; MLP, multi-layer perceptron.

to assess the output sequence and provide feedback to refine the model. They examined the model for six folds using Rosetta-based structure prediction, and gcWGAN had higher TM score distributions and more diverse sequence profiles than CVAE.<sup>195</sup> Another notable experiment is Ingraham et al.'s<sup>23</sup> graph transformer model that inputs a structure, represented as a graph, and outputs the sequence profile. They treat the sequence design problem similar to a machine translation problem, i.e., a translation from structure to sequence. Like the original transformer model,<sup>57</sup> they adopted an encoder-decoder framework with self-attention mechanisms to dynamically learn the relationship between information in two neighbor layers. They measured their results by perplexity, a widely used metric in speech recognition,<sup>212</sup> and the per-residue perplexity (lower is better) for single chains was 9.15, lower than the perplexity for SPIN2 (12.86). Norn et al. treated the protein design problem as that of maximizing the probability of a sequence given a structure. They back-propagate through the trRosetta structure prediction network<sup>103</sup> to find a sequence that minimizes the distance between predicted structure and a desired structure.<sup>203</sup> Norn et al. validate their designs computationally by showing the generated sequences have deep wells in their modeled energy landscapes. Strokach et al. treated the design of protein sequence given a target structure as a constraint satisfaction problem. They optimized their GNN architecture on the related problem of filling in a Sudoku puzzle followed by training on millions of protein sequences corresponding to thousands of structural folds. They were able to validate designed sequences *in silico* and demonstrate that some designs folded to their target structures *in vitro*.<sup>213</sup>

An ambitious design goal is to generate new structures without specifying the target structure. Anand and Huang were the first to generate new structures using DL. They tested various representations (e.g., full atom, torsion-only) with a deep convolutional GAN (DCGAN) framework that generates sequenceagnostic, fixed-length short protein structural fragments.<sup>24</sup> They found that the distance map of  $C_{\alpha}$  atoms gives the most meaningful protein structures, although the asymmetry of  $\psi$ and  $\phi$  torsion angles<sup>214</sup> was only recovered with torsion-based representations. Later they extended this work to all atoms in the backbone and combined with a recovery network to avoid the time-consuming structure reconstruction process.<sup>68</sup> They showed that some of the designed folds are stable in molecular simulation. In a more narrowly focused study, Eguchi et al.<sup>192</sup> trained a VAE model with the structures of immunoglobulin (Ig) proteins, called Ig-VAE. By sampling the latent space, they



generated 5,000 new Ig structures (sequence-agnostic) and then screened them with computational docking to identify putative binders to SARS-CoV2-RBD.

Another approach exploits a DL structure prediction algorithm and a Markov Chain MC (MCMC) search to find sequences that fold into novel compact structures. Anishchenko et al.<sup>193</sup> iterated sequences through the DL network, trRosetta,<sup>103</sup> to "hallucinate"<sup>215</sup> mutually compatible sequence-structure pairs in a manner similar to "input design".<sup>183</sup> By maximizing the contrast between the distance distributions predicted by trRosetta and a background network trained on noise, they obtained new sequences with geometric maps with sharp geometric features. Impressively, 27 of the<sup>127</sup> hallucinated sequences were experimentally validated to fold into monomeric, highly stable, proteins with circular dichroism spectra compatible with the predicted structure.

#### **OUTLOOK AND CONCLUSION**

In this review, we have summarized the current state-of-the-art DL techniques applied to the problem of protein structure prediction and design. As in many other areas, DL shows the potential to revolutionize the field of protein modeling. While DL originated from computer vision, NLP and machine learning, its fast development combined with knowledge from operations research,<sup>216</sup> game theory,<sup>65</sup> and variational inference<sup>32</sup> among other fields, has resulted in many new and powerful frameworks to solve increasingly complex problems. The application of DL for biomolecular structure has just begun, and we expect to see more efforts on methodology development and applications in protein modeling and design.

We observed several trends.

#### **Experimental Validation**

An important gap in current DL work in protein modeling, especially protein design (with few notable exceptions),<sup>205,190,198,193</sup> is the lack of experimental validation. Past blind challenges, e.g., CASP and CAPRI, and design claims have shown that experimental validation in this field is of paramount importance, where computational models are still prone to error. A key next stage for this field is to engage collaborations between machine learning experts and experimental protein engineers to test and validate these emerging approaches.

#### Importance of Benchmarking

In other fields of machine learning, standardized benchmarks have triggered rapid progress.<sup>217–219</sup> CASP is a great example that provides a standardized platform for benchmarking diverse algorithms, including emerging DL-based approaches. A welldefined question and proper evaluation (especially experimental) would lead to more open competition among a broader range of groups and, eventually, the innovation of more diverse and powerful algorithms.

#### **Imposing a Physics-Based Prior**

One common topic among the machine learning community is how to utilize existing domain knowledge to reduce the effort during training. Unlike certain classical ML problems, such as image classification, in protein modeling, a wide range of biophys-



ical principles restrict the range of plausible solutions. Some examples in related fields include imposing a physics-based model prior,<sup>220,221</sup> adding a regularization term with physical meaning,<sup>222</sup> and adopting a specific formula to conserve physical symmetry.<sup>223,224</sup> Similarly, in protein modeling, well-established empirical observations can help restrict the solution space, such as the Ramanchandran distribution of backbone torsion angles<sup>214</sup> and the Dunbrack or Richardsons library of side-chain conformations.<sup>225,226</sup>

#### **Closed-Loop Design**

The performance of DL methodologies relies heavily on the quality of data, but the publicly available datasets may not cover important sample space because of experimental accessibility at the time of experiments. Furthermore, the dataset may contain harmful noise from non-uniform experimental protocols and conditions. A possible solution may be to combine model training with experimental data generation. For instance, one may devise a closed-loop strategy to generate experimental data, on-the-fly, for gueries (or model inputs) that are most likely to improve the model, and update the training dataset with the newly generated data.<sup>227-230</sup> For such a strategy to be feasible, automated synthesis and characterization is necessary. As highthroughput synthesis and testing of protein (or DNA and RNA) can be carried out in parallel, automation is possible. While such a strategy may seem far-fetched, automated platforms such as those from Ginkgo Bioworks or Transcriptic are already on the market.

#### **Reinforcement Learning**

Another approach to overcome the limitation of data availability is reinforcement learning (RL). Biologically meaningful data may be generated on-the-fly in simulated environments, such as the Foldit game. In the most famous application of RL, AlphaGo Zero.<sup>21</sup> an RL agent (network) was able to learn and master the game by learning from the game environment alone. There are already some examples of RL in the field of chemistry and electric engineering to optimize organic molecules or computational chips.<sup>231-233</sup> One suitable protein modeling problem for an RL algorithm would be training an artificial intelligence (AI) agent to make a series of "moves" to fold a protein, similar to the Foldit game.<sup>234,235</sup> Such studies are still rare and previous attempts have focused on folding the 2D hydrophobic-polar model of proteins.<sup>236,237</sup> Although the results did not yet beat conventional methods, Gao<sup>238</sup> recently explored using policy and reward networks in an RL scheme to fold 3D protein structures de novo by guiding the selection of MC moves in Rosetta. Angermueller et al.239 applied a model-based RL framework to designing sequences of AMPs and transcription factor binding sites.

#### Model Interpretability

One should keep in mind that a neural network represents nothing more (and nothing less) than a powerful and flexible regression model. In addition, due to their highly recursive nature, neural networks tend to be regarded as "black-boxes", i.e., too complicated for practitioners to understand the resulting parameters and functions. Although model interpretability in ML





is a rapidly developing field, many popular approaches, such as saliency analysis<sup>240–242</sup> for image classification models, are far from satisfactory.<sup>243</sup> Although other approaches<sup>244,245</sup> offer more reliable interpretations, their application to DL model interpretation has been largely missing in protein modeling. As a result, current DL models offer limited understanding of the complex patterns they learn.

#### **Beyond Proteins**

DL-based methods are general and so, with appropriate representation and sufficient training data, they can be applied to other molecules. Like proteins, nucleic acids, carbohydrates, and lipids are also polymers, composed of nucleotides, monosaccharides, and aliphatic subunits and head groups, respectively. Many approaches developed for learning protein sequence and structural information can be extended to these other classes of biomolecules.<sup>246,247</sup> Finally, biology often conjugates these molecules, e.g., for glycoproteins. DL approaches that build up from basic chemistry, such as those being developed for small molecule drugs,<sup>248–251</sup> may inspire methods to treat these biomolecules that do not fall into a strict polymer type.

#### The "Sequence $\rightarrow$ Structure $\rightarrow$ Function" Paradigm

We know from molecular biophysics that a sequence translates into function through the physical intermediary of a 3D molecular structure. Allosteric proteins,252 for instance, may exhibit different structural conformations under different physiological conditions (e.g., pH) or environmental stimuli (e.g., small molecules, inhibitors), reminding us that context is as important as protein sequence. That is, despite Anfinsen's<sup>42</sup> hypothesis, sequence alone does not always fully determine the structure. Some proteins require chaperones to fold to their native structure, meaning that a sequence could result in non-native conformations when the kinetics of folding to the native structure may be unfavorable in the absence of a chaperone. Because many powerful DL algorithms in NLP operate on sequential data, it may seem reasonable to use protein sequences alone for training DL models. In principle, with a suitable framework and training, DL could disentangle the underlying relationships between sequence and structural elements. However, a careful selection of DL frameworks that are structure or mechanism-aware will accelerate learning and improve predictive power. Indeed, many successful DL frameworks applied so far (e.g., CNNs or graph CNNs) factor in the importance of learning on structural information.

Finally, with the hope of gaining insight into the fundamental science of biomolecules, there is a desire to link AI approaches to the underlying biochemical and biophysical principles that drive biomolecular function. For more practical purposes, a deeper understanding of underlying principles and hidden patterns that lead to pathology is important in the development of therapeutics. Thus, while efforts strictly limited to sequences are abundant, we believe that models with structural insights will play a more critical role in the future.

#### ACKNOWLEDGMENTS

This work was supported by the NIH through grant R01-GM078221. We thank Dr. Justin S. Smith at the Center for Nonlinear Studies at Los Alamos National Laboratory, NM, for helpful discussion and Dr. Andrew D. White at the Department of Chemical Engineering at University of Rochester, NY, and Alexander Rives at the Department of Computer Science at New York University, NY, for helpful suggestions. We are also grateful for insightful suggestions from the reviewers.

#### **AUTHOR CONTRIBUTIONS**

Conceptualization, W.G. and J.J.G.; Investigation, W.G. and S.P.M.; Writing – Original Draft, W.G.; Writing – Review & Editing, W.G., S.P.M., J.S., and J.J.G.; Funding Acquisition, J.J.G.; Resources, J.J.G.; Supervision, J.S. and J.J.G.

#### REFERENCES

- Slabinski, L., Jaroszewski, L., Rodrigues, A.P., Rychlewski, L., Wilson, I.A., Lesley, S.A., and Godzik, A. (2007). The challenge of protein structure determination-lessons from structural genomics. Protein Sci. 16, 2472–2482.
- Markwick, P.R.L., Malliavin, T., and Nilges, M. (2008). Structural biology by NMR: structure, dynamics, and interactions. PLoS Comput. Biol. 4, e1000168.
- Jonic, S., and Vénien-Bryan, C. (2009). Protein structure determination by electron cryo-microscopy. Curr. Opin. Pharmacol. 9, 636–642.
- Kryshtafovych, A., Schwede, T., Topf, M., Fidelis, K., and Moult, J. (2019). Critical assessment of methods of protein structure prediction (CASP)—Round XIII. Proteins 87, 1011–1020.
- Hollingsworth, S.A., and Dror, R.O. (2018). Molecular dynamics simulation for all. Neuron 99, 1129–1143.
- Ranjan, A., Fahad, M.S., Fernandez-Baca, D., Deepak, A., and Tripathi, S. (2019). Deep robust framework for protein function prediction using variable-length protein sequences. IEEE/ACM Trans. Comput. Biol. Bioinform. 17, 1648–1659.
- 7. Huang, P.S., Boyken, S.E., and Baker, D. (2016). The coming of age of de novo protein design. Nature 537, 320–327.
- Yang, K.K., Wu, Z., and Arnold, F.H. (2019). Machine-learning-guided directed evolution for protein engineering. Nat. Methods 16, 687–694.
- Bohr, H., Bohr, J., Brunak, S., Cotterill, J.R., Fredholm, H., Lautrup, B., and Petersen, S. (1990). A novel approach to prediction of the 3-dimensional structures of protein backbones by neural networks. FEBS Lett. 261, 43–46.
- Schneider, G., and Wrede, P. (1994). The rational design of amino acid sequences by artificial neural networks and simulated molecular evolution: de novo design of an idealized leader peptidase cleavage site. Biophys. J. 66, 335–344.
- Schneider, G., Schrödl, W., Wallukat, G., Müller, J., Nissen, E., Rönspeck, W., Wrede, P., and Kunze, R. (1998). Peptide design by artificial neural networks and computer-based evolutionary search. Proc. Natl. Acad. Sci. U S A *95*, 12179–12184.
- Ofran, Y., and Rost, B. (2003). Predicted protein-protein interaction sites from local sequence information. FEBS Lett. 544, 236–239.
- Nielsen, M., Lundegaard, C., Worning, P., Lauemøller, S.L., Lamberth, K., Buus, S., Brunak, S., and Lund, O. (2003). Reliable prediction of T-cell epitopes using neural networks with novel sequence representations. Protein Sci. 12, 1007–1017.
- 14. LeCun, Y., Bengio, Y., and Hinton, G. (2015). Deep learning. Nature 521, 436–444.
- 15. Angermueller, C., Pärnamaa, T., Parts, L., and Stegle, O. (2016). Deep learning for computational biology. Mol. Syst. Biol. *12*, 878.
- Ching, T., Himmelstein, D.S., Beaulieu-Jones, B.K., Kalinin, A.A., Do, B.T., Way, G.P., Ferrero, E., Agapow, P.-M., Zietz, M., Hoffman, M.M., et al. (2018). Opportunities and obstacles for deep learning in biology and medicine. J. R. Soc. Interfaces 15, 20170387.





- Mura, C., Draizen, E.J., and Bourne, P.E. (2018). Structural biology meets data science: does anything change? Curr. Opin. Struct. Biol. 52, 95–102.
- Noé, F., De Fabritiis, G., and Clementi, C. (2020). Machine learning for protein folding and dynamics. Curr. Opin. Struct. Biol. 60, 77–84.
- 19. Guo, Y., Liu, Y., Oerlemans, A., Lao, S., Wu, S., and Lew, M.S. (2016). Deep learning for visual understanding: a review. Neurocomputing 187, 27–48.
- Young, T., Hazarika, D., Poria, S., and Cambria, E. (2018). Recent trends in deep learning based natural language processing. IEEE Comput. Intelligence Mag. 13, 55–75.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. (2017). Mastering the game of go without human knowledge. Nature 1550, 354.
- 22. Senior, A.W., Evans, R., Jumper, J., Kirkpatrick, J., Sifre, L., Green, T., Chongli, Q., Zidek, A., Nelson, A.W.R., Bridgland, A., et al. (2019). Protein structure prediction using multiple deep neural networks in the 13th Critical Assessment of Protein Structure Prediction (CASP13). Proteins 87, 1141–1148.
- Ingraham, J., Garg, V., Barzilay, R., and Jaakkola, T. (2019). Generative models for graph-based protein design. Adv. Neural Inf. Process. Syst. 15820–15831.
- Anand, N., and Huang, P. (2018). Generative modeling for protein structures. Adv. Neural Inf. Process. Syst. 7494–7505.
- 25. O'Connell, J., Li, Z., Hanson, J., Heffernan, R., Lyons, J., Paliwal, K., Dehzangi, A., Yang, Y., and Zhou, Y. (2018). SPIN2: predicting sequence profiles from protein structures using deep neural networks. Proteins: Struct. Funct. Bioinformatics 86, 629–633.
- Senior, A.W., Evans, R., Jumper, J., Kirkpatrick, J., Sifre, L., Green, T., Qin, C., Žídek, A., Nelson, A.W., Bridgland, A., et al. (2020). Improved protein structure prediction using potentials from deep learning. Nature, 1–5.
- 27. Li, Y., Huang, C., Ding, L., Li, Z., Pan, Y., and Gao, X. (2019). Deep learning in bioinformatics: introduction, application, and perspective in the big data era. Methods *166*, *4*–21.
- Noé, F., Tkatchenko, A., Müller, K.-R., and Clementi, C. (2020). Machine learning for molecular simulation. Annu. Rev. Phys. Chem. 71, 361–390.
- Graves, J., Byerly, J., Priego, E., Makkapati, N., Parish, S.V., Medellin, B., and Berrondo, M. (2020). A review of deep learning methods for antibodies. Antibodies 9, 12.
- Kandathil, S.M., Greener, J.G., and Jones, D.T. (2019). Recent developments in deep learning applied to protein structure prediction. Proteins: Struct. Funct. Bioinformatics 87, 1179–1189.
- Torrisi, M., Pollastri, G., and Le, Q. (2020). Deep learning methods in protein structure prediction. Comput. Struct. Biotechnol. J. 18, 1301–1310.
- Kingma, D.P., and Welling, M. (2013). Auto-encoding variational Bayes. arXiv 1312, 6114.
- **33.** Pauling, L., and Niemann, C. (1939). The structure of proteins. J. Am. Chem. Soc. *61*, 1860–1867.
- Kuhlman, B., and Bradley, P. (2019). Advances in protein structure prediction and design. Nat. Rev. Mol. Cell Biol. 20, 681–697.
- UniProt-Consortium. (2019). UniProt: a worldwide hub of protein knowledge. Nucleic Acids Res. 47, D506–D515.
- Kuhlman, B., Dantas, G., Ireton, G.C., Varani, G., Stoddard, B.L., and Baker, D. (2003). Design of a novel globular protein fold with atomic-level accuracy. Science 302, 1364–1368.
- Fisher, M.A., McKinley, K.L., Bradley, L.H., Viola, S.R., and Hecht, M.H. (2011). De novo designed proteins from a library of artificial sequences function in Escherichia coli and enable cell growth. PLoS One 6, e15364.
- Correia, B.E., Bates, J.T., Loomis, R.J., Baneyx, G., Carrico, C., Jardine, J.G., Rupert, P., Correnti, C., Kalyuzhniy, O., Vittal, V., et al. (2014). Proof of principle for epitope-focused vaccine design. Nature 507, 201.

- 39. King, N.P., Sheffler, W., Sawaya, M.R., Vollmar, B.S., Sumida, J.P., André, I., Gonen, T., Yeates, T.O., and Baker, D. (2012). Computational design of self-assembling protein nanomaterials with atomic level accuracy. Science 336, 1171–1174.
- 40. Tinberg, C.E., Khare, S.D., Dou, J., Doyle, L., Nelson, J.W., Schena, A., Jankowski, W., Kalodimos, C.G., Johnsson, K., Stoddard, B.L., et al. (2013). Computational design of ligand-binding proteins with high affinity and selectivity. Nature 501, 212–216.
- Joh, N.H., Wang, T., Bhate, M.P., Acharya, R., Wu, Y., Grabe, M., Hong, M., Grigoryan, G., and DeGrado, W.F. (2014). De novo design of a transmembrane Zn<sup>2+</sup>-transporting four-helix bundle. Science 346, 1520–1524.
- 42. Anfinsen, C.B. (1973). Principles that govern the folding of protein chains. Science 181, 223–230.
- Levinthal, C. (1968). Are there pathways for protein folding? J. Chim. Phys. 65, 44–45.
- 44. Li, B., Fooksa, M., Heinze, S., and Meiler, J. (2018). Finding the needle in the haystack: towards solving the protein-folding problem computationally. Crit. Rev. Biochem. Mol. Biol. 53, 1–28.
- Dahiyat, B.I., and Mayo, S.L. (1997). De novo protein design: fully automated sequence selection. Science 278, 82–87.
- Korendovych, I.V., and DeGrado, W.F. (2020). De novo protein design, a retrospective. Q. Rev. Biophys. 53, https://doi.org/10.1017/ S0033583519000131.
- Dougherty, M.J., and Arnold, F.H. (2009). Directed evolution: new parts and optimized function. Curr. Opin. Biotechnol. 20, 486–491.
- Sun, R. (2019). Optimization for deep learning: theory and algorithms. ar-Xiv 1912, 08957.
- 49. Schmidhuber, J. (2015). Deep learning in neural networks: an overview. Neural Networks 61, 85–117.
- LeCun, Y., Boser, B.E., Denker, J.S., Henderson, D., Howard, R.E., Hubbard, W.E., and Jackel, L.D. (1990). Handwritten digit recognition with a back-propagation network. Adv. Neural Inf. Process. Syst. 396–404.
- He, K.; Zhang, X.; Ren, S.; Sun, J. Deep residual learning for image recognition. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition 2016, 770–778.
- Jordan, M.I. (1997). Serial order: a parallel distributed processing approach. Adv. Psychol. 121, 471–495.
- Hochreiter, S., and Schmidhuber, J. (1997). Long short-term memory. Neural Comput. 9, 1735–1780.
- Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. (2014). Learning phrase representations using RNN encoder-decoder for statistical machine translation. arXiv 1406, 1078.
- Müller, A.T., Hiss, J.A., and Schneider, G. (2018). Recurrent neural network model for constructive peptide design. J. Chem. Inf. Model. 58, 472–479.
- Bahdanau, D.; Cho, K.H.; Bengio, Y. Neural machine translation by jointly learning to align and translate. 3rd International Conference on Learning Representations, ICLR 2015—Conference Track Proceedings. 2015.
- 57. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. Adv. Neural Inf. Process. Syst. 2017, 5999–6009.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). Bert: pretraining of deep bidirectional transformers for language understanding. arXiv 1810, 04805.
- Rives, A., Goyal, S., Meier, J., Guo, D., Ott, M., Zitnick, C.L., et al. (2019). Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. bioRxiv, 622803, https://doi. org/10.1101/622803. https://www.biorxiv.org/content/10.1101/622803v3.

# Patterns

Review

- **60.** Pittala, S., and Bailey-Kellogg, C. (2020). Learning context-aware structural representations to predict antigen and antibody binding interfaces. Bioinformatics *36*, 3996–4003.
- Hinton, G.E., and Zemel, R.S. (1994). Autoencoders, minimum description length and Helmholtz free energy. Adv. Neural Inf. Process. Syst. 3–10.
- Kingma, D.P., and Welling, M. (2019). An introduction to variational autoencoders. arXiv 1906, 02691.
- Blei, D.M., Kucukelbir, A., and McAuliffe, J.D. (2017). Variational inference: a review for statisticians. J. Am. Stat. Assoc. 112, 859–877.
- Das, P., Wadhawan, K., Chang, O., Sercu, T., Santos, C.D., Riemer, M., Chenthamarakshan, V., Padhi, I., and Mojsilovic, A. (2018). PepCVAE: semi-supervised targeted design of antimicrobial peptide sequences. ar-Xiv 1810, 07743.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., et al. (2014). Generative adversarial nets. Adv. Neural Inf. Process. Syst. 2672–2680. https://papers.nips.cc/paper/5423-generativeadversarial-nets.
- Arjovsky, M., Chintala, S., Bottou, L., and Wasserstein, G.A.N. (2017). ar-Xiv 1701, 07875.
- Kurach, K., Lučić, M., Zhai, X., Michalski, M., and Gelly, S. (2019). A largescale study on regularization and normalization in GANs. Int. Conf. Mach. Learn. 3581–3590.
- Anand, N., Eguchi, R., and Huang, P.-S. (2019). Fully differentiable fullatom protein backbone generation. Int. Conf. Learn. Rep. 35. https:// openreview.net/revisions?id=SJxnVL8YOV.
- Niepert, M., Ahmed, M., and Kutzkov, K. (2016). Learning convolutional neural networks for graphs. Int. Conf. Mach. Learn. 2014–2023.
- Luo, F., Wang, M., Liu, Y., Zhao, X.-M., and Li, A. (2019). DeepPhos: prediction of protein phosphorylation sites with deep learning. Bioinformatics 35, 2766–2773.
- 71. Li, F., Chen, J., Leier, A., Marquez-Lago, T., Liu, Q., Wang, Y., Revote, J., Smith, A.I., Akutsu, T., Webb, G.I., et al. (2020). DeepCleave: a deep learning predictor for caspase and matrix metalloprotease substrates and cleavage sites. Bioinformatics 36, 1057–1065.
- Bengio, Y., Courville, A., and Vincent, P. (2013). Representation learning: a review and new perspectives. IEEE Trans. Pattern Anal. Mach. Intelligence 35, 1798–1828.
- Romero, P.A., Krause, A., and Arnold, F.H. (2013). Navigating the protein fitness landscape with Gaussian processes. Proc. Natl. Acad. Sci. U S A 110, E193–E201.
- Bedbrook, C.N., Yang, K.K., Rice, A.J., Gradinaru, V., and Arnold, F.H. (2017). Machine learning to design integral membrane channel rhodopsins for efficient eukaryotic expression and plasma membrane localization. PLoS Comput. Biol. *13*, e1005786.
- Ofer, D., and Linial, M. (2015). ProFET: feature engineering captures highlevel protein functions. Bioinformatics 31, 3429–3436.
- Kawashima, S., Pokarowski, P., Pokarowska, M., Kolinski, A., Katayama, T., and Kanehisa, M. (2007). AAindex: amino acid index database, progress report 2008. Nucleic Acids Res. 36, D202–D205.
- Wang, S., Peng, J., Ma, J., and Xu, J. (2016). Protein secondary structure prediction using deep convolutional neural fields. Sci. Rep. 6, 18962.
- Drori, I., Thaker, D., Srivatsa, A., Jeong, D., Wang, Y., Nan, L., Wu, F., Leggas, D., Lei, J., Lu, W., et al. (2019). Accurate protein structure prediction by embeddings and deep learning representations. arXiv 1911, 05531.
- 79. Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient estimation of word representations in vector space. arXiv 1301, 3781.
- Le, Q., and Mikolov, T. (2014). Distributed representations of sentences and documents. Int. Conf. Mach. Learn. 1188–1196.



- Asgari, E., and Mofrad, M.R. (2015). Continuous distributed representation of biological sequences for deep proteomics and genomics. PLoS One 10, e0141287.
- 82. El-Gebali, S., Mistry, J., Bateman, A., Eddy, S.R., Luciani, A., Potter, S.C., Qureshi, M., Richardson, L.J., Salazar, G.A., Smart, A., et al. (2019). The Pfam protein families database in 2019. Nucleic Acids Res. 47, D427–D432.
- Cai, C., Han, L., Ji, Z.L., Chen, X., and Chen, Y.Z. (2003). SVM-Prot: webbased support vector machine software for functional classification of a protein from its primary sequence. Nucleic Acids Res. 31, 3692–3697.
- Aragues, R., Sali, A., Bonet, J., Marti-Renom, M.A., and Oliva, B. (2007). Characterization of protein hubs by inferring interacting motifs from protein interactions. PLoS Comput. Biol. 3, e178.
- Yu, C., van der Schaar, M., and Sayed, A.H. (2016). Distributed learning for stochastic generalized Nash equilibrium problems. CoRR. https://doi. org/10.1109/TSP.2017.2695451.
- Yang, K.K., Wu, Z., Bedbrook, C.N., and Arnold, F.H. (2018). Learned protein embeddings for machine learning. Bioinformatics 34, 2642–2648.
- Alley, E.C., Khimulya, G., Biswas, S., AlQuraishi, M., and Church, G.M. (2019). Unified rational protein engineering with sequence-based deep representation learning. Nat. Methods 16, 1315–1322.
- Krause, B., Lu, L., Murray, I., and Renals, S. (2016). Multiplicative LSTM for sequence modelling. arXiv 1609, 07959.
- Heinzinger, M., Elnaggar, A., Wang, Y., Dallago, C., Nechaev, D., Matthes, F., and Rost, B. (2019). Modeling aspects of the language of life through transfer-learning protein sequences. BMC Bioinformatics 20, 723.
- Peters, M.E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., and Zettlemoyer, L. (2018). Deep contextualized word representations. arXiv 1802, 05365.
- Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. arXiv 2005, 14165.
- Ding, X., Zou, Z., and Brooks, C.L., lii (2019). Deciphering protein evolution and fitness landscapes with latent space models. Nat. Commun. 210, 1–13.
- Sinai, S., Kelsic, E., Church, G.M., and Nowak, M.A. (2017). Variational auto-encoding of protein sequences. arXiv 1712, 03346.
- Riesselman, A.J., Ingraham, J.B., and Marks, D.S. (2018). Deep generative models of genetic variation capture the effects of mutations. Nat. Methods 15, 816–822.
- Rao, R., Bhattacharya, N., Thomas, N., Duan, Y., Chen, P., Canny, J., et al. (2019). Evaluating protein transfer learning with TAPE. Adv. Neural Inf. Process. Syst. 9689–9701. http://papers.nips.cc/paper/9163evaluating-protein-transfer-learning-with-tape.
- 96. Townshend, R., Bedi, R., and Dror, R.O. (2018). Generalizable protein interface prediction with end-to-end learning. arXiv 1807, 01297.
- Simonovsky, M., and Meyers, J. (2020). DeeplyTough: learning structural comparison of protein binding sites. J. Chem. Inf. Model. 60, 2356–2366.
- Kolodny, R., Koehl, P., Guibas, L., and Levitt, M. (2002). Small libraries of protein fragments model native protein structures accurately. J. Mol. Biol. 323, 297–307.
- 99. Taylor, W.R.A. (2002). "periodic table" for protein structures. Nature 416, 657–660.
- 100. Li, J., and Koehl, P. (2014). 3D representations of amino acidsapplications to protein sequence comparison and classification. Comput. Struct. Biotechnol. J. 11, 47–58.
- AlQuraishi, M. (2019). End-to-End differentiable learning of protein structure. Cell Syst. 8, 292–301.e3.
- 102. Wang, S., Sun, S., Li, Z., Zhang, R., and Xu, J. (2017). Accurate de novo prediction of protein contact map by ultra-deep learning model. PLoS Comput. Biol. 13, e1005324.





- 103. Yang, J., Anishchenko, I., Park, H., Peng, Z., Ovchinnikov, S., and Baker, D. (2020). Improved protein structure prediction using predicted interresidue orientations. Proc. Natl. Acad. Sci. U S A 117, 1496–1503.
- Brunger, A.T. (2007). Version 1.2 of the crystallography and NMR system. Nat. Protoc. 2, 2728.
- 105. Zhou, J., Cui, G., Zhang, Z., Yang, C., Liu, Z., and Sun, M. (2018). Graph neural networks: a review of methods and applications. arXiv 1812, 08434.
- 106. Ahmed, E., Saint, A., Shabayek, A., Cherenkova, K., Das, R., Gusev, G., Aouada, D., and Ottersten, B. (2018). Deep learning advances on different 3D data representations: a survey. arXiv 1, 01462.
- Wu, Z., Pan, S., Chen, F., Long, G., Zhang, C., and Philip, S.Y. (2020). A comprehensive survey on graph neural networks. IEEE Trans. Neural Networks Learn. Syst. 1–21. https://ieeexplore.ieee.org/abstract/ document/9046288.
- Vishveshwara, S., Brinda, K., and Kannan, N. (2002). Protein structure: insights from graph theory. J. Theor. Comput. Chem. 1, 187–211.
- 109. Ying, Z., You, J., Morris, C., Ren, X., Hamilton, W., and Leskovec, J. (2018). Hierarchical graph representation learning with differentiable pooling. Adv. Neural Inf. Process. Syst. 4800–4810. https://papers. nips.cc/paper/7729-hierarchical-graph-representation-learning-withdifferentiable-pooling.
- Borgwardt, K.M., Ong, C.S., Schönauer, S., Vishwanathan, S., Smola, A.J., and Kriegel, H.-P. (2005). Protein function prediction via graph kernels. Bioinformatics 21, i47–i56.
- 111. Dobson, P.D., and Doig, A.J. (2003). Distinguishing enzyme structures from non-enzymes without alignments. J. Mol. Biol. 330, 771–783.
- 112. Fout, A., Byrd, J., Shariat, B., and Ben-Hur, A. (2017). Protein interface prediction using graph convolutional networks. Adv. Neural Inf. Process. Syst. 6530–6539. https://papers.nips.cc/paper/7231-protein-interfaceprediction-using-graph-convolutional-networks.
- 113. Zamora-Resendiz, R., and Crivelli, S. (2019). Structural learning of proteins using graph convolutional neural networks. bioRxiv, 610444. https://www.biorxiv.org/content/10.1101/610444v1.
- Gligorijevic, V., Renfrew, P.D., Kosciolek, T., Leman, J.K., Cho, K., Vatanen, T., et al. (2019). Structure-based function prediction using graph convolutional networks. bioRxiv, 786236. https://www.biorxiv.org/ content/10.1101/786236v2.
- Torng, W., and Altman, R.B. (2019). Graph convolutional neural networks for predicting drug-target interactions. J. Chem. Inf. Model. 59, 4131–4149.
- 116. Gainza, P., Sverrisson, F., Monti, F., Rodola, E., Boscaini, D., Bronstein, M., and Correia, B. (2020). Deciphering interaction fingerprints from protein molecular surfaces using geometric deep learning. Nat. Methods 17, 184–192.
- Bronstein, M.M., Bruna, J., LeCun, Y., Szlam, A., and Vandergheynst, P. (2017). Geometric deep learning: going beyond Euclidean data. IEEE Signal. Process. Mag. 34, 18–42.
- Nerenberg, P.S., and Head-Gordon, T. (2018). New developments in force fields for biomolecular simulations. Curr. Opin. Struct. Biol. 49, 129–138.
- Derevyanko, G., Grudinin, S., Bengio, Y., and Lamoureux, G. (2018). Deep convolutional networks for quality assessment of protein folds. Bioinformatics 34, 4046–4053.
- **120.** Best, R.B., Zhu, X., Shim, J., Lopes, P.E., Mittal, J., Feig, M., and MacKerell, A.D., Jr. (2012). Optimization of the additive CHARMM all-atom protein force field targeting improved sampling of the backbone  $\phi$ ,  $\psi$  and side-chain  $\chi 1$  and  $\chi 2$  dihedral angles. J. Chem. Theor. Comput. 8, 3257–3273.
- 121. Weiner, S.J., Kollman, P.A., Case, D.A., Singh, U.C., Ghio, C., Alagona, G., Profeta, S., and Weiner, P. (1984). A new force field for molecular mechanical simulation of nucleic acids and proteins. J. Am. Chem. Soc. 106, 765–784.

- 122. Alford, R.F., Leaver-Fay, A., Jeliazkov, J.R., O'Meara, M.J., DiMaio, F.P., Park, H., Shapovalov, M.V., Renfrew, P.D., Mulligan, V.K., Kappel, K., et al. (2017). The Rosetta all-atom energy function for macromolecular modeling and design. J. Chem. Theor. Comput. *13*, 3031–3048.
- 123. Behler, J., and Parrinello, M. (2007). Generalized neural-network representation of high-dimensional potential-energy surfaces. Phys. Rev. Lett. 98, 146401.
- 124. Smith, J.S., Isayev, O., and Roitberg, A.E. (2017). ANI-1: an extensible neural network potential with DFT accuracy at force field computational cost. Chem. Sci. 8, 3192–3203.
- 125. Smith, J.S., Nebgen, B., Lubbers, N., Isayev, O., and Roitberg, A.E. (2018). Less is more: sampling chemical space with active learning. J. Chem. Phys. 148, 241733.
- 126. Schütt, K.T., Arbabzadah, F., Chmiela, S., Müller, K.R., and Tkatchenko, A. (2017). Quantum-chemical insights from deep tensor neural networks. Nat. Commun. 8, 1–8.
- 127. Schütt, K.T., Sauceda, H.E., Kindermans, P.-J., Tkatchenko, A., and Müller, K.-R. (2018). SchNet—a deep learning architecture for molecules and materials. J. Chem. Phys. *148*, 241722.
- 128. Zhang, L., Han, J., Wang, H., Car, R., and Weinan, E. (2018). Deep potential molecular dynamics: a scalable model with the accuracy of quantum mechanics. Phys. Rev. Lett. *120*, 143001.
- Unke, O.T., and Meuwly, M. (2019). PhysNet: a neural network for predicting energies, forces, dipole moments, and partial charges. J. Chem. Theor. Comput. 15, 3678–3693.
- 130. Zubatyuk, R., Smith, J.S., Leszczynski, J., and Isayev, O. (2019). Accurate and transferable multitask prediction of chemical properties with an atoms-in-molecules neural network. Sci. Adv. 5, eaav6490.
- Lahey, S.-L.J., and Rowley, C.N. (2020). Simulating protein-ligand binding with neural network potentials. Chem. Sci. 11, 2362–2368.
- 132. Wang, Z., Han, Y., Li, J., and He, X. (2020). Combining the fragmentation approach and neural network potential energy surfaces of fragments for accurate calculation of protein energy. J. Phys. Chem. B 124, 3027–3035.
- 133. Senn, H.M., and Thiel, W. (2009). QM/MM methods for biomolecular systems. Angew. Chem. Int. Ed. 48, 1198–1229.
- Wang, Y., Fass, J., and Chodera, J.D. (2020). End-to-End Differentiable Molecular Mechanics Force Field Construction (arXiv). https://arxiv.org/ abs/2010.01196.
- 135. Kmiecik, S., Gront, D., Kolinski, M., Wieteska, L., Dawid, A.E., and Kolinski, A. (2016). Coarse-grained protein models and their applications. Chem. Rev. 116, 7898–7936.
- **136.** Zhang, L., Han, J., Wang, H., Car, R., and Weinan, E. (2018). DeePCG: constructing coarse-grained models via deep neural networks. J. Chem. Phys. *149*, 034101.
- 137. Patra, T.K., Loeffler, T.D., Chan, H., Cherukara, M.J., Narayanan, B., and Sankaranarayanan, S.K. (2019). A coarse-grained deep neural network model for liquid water. Appl. Phys. Lett. *115*, 193101.
- Wang, J., Olsson, S., Wehmeyer, C., Pérez, A., Charron, N.E., De Fabritiis, G., Noé, F., and Clementi, C. (2019). Machine learning of coarsegrained molecular dynamics force fields. ACS Cent. Sci. 5, 755–767.
- 139. Wang, W., and Gómez-Bombarelli, R. (2019). Learning coarse-grained particle latent space with auto-encoders. Adv. Neural Inf. Process. Syst. 1.
- 140. Li, Z., Wellawatte, G.P., Chakraborty, M., Gandhi, H.A., Xu, C., and White, A.D. (2020). Graph neural network based coarse-grained mapping prediction. Chem. Sci. 11, 9524–9531.
- 141. Jones, D.T., Buchan, D.W., Cozzetto, D., and Pontil, M. (2011). PSICOV: precise structural contact prediction using sparse inverse covariance estimation on large multiple sequence alignments. Bioinformatics 28, 184–190.
- 142. Di Lena, P., Nagata, K., and Baldi, P. (2012). Deep architectures for protein contact map prediction. Bioinformatics 28, 2449–2457.

# Patterns

Review

- 143. Eickholt, J., and Cheng, J. (2012). Predicting protein residue-residue contacts using deep networks and boosting. Bioinformatics 28, 3066–3072.
- 144. Seemayer, S., Gruber, M., and Söding, J. (2014). CCMpred—fast and precise prediction of protein residue-residue contacts from correlated mutations. Bioinformatics 30, 3128–3130.
- 145. Skwark, M.J., Raimondi, D., Michel, M., and Elofsson, A. (2014). Improved contact predictions using the recognition of protein like contact patterns. PLoS Comput. Biol. 10, e1003889.
- 146. Jones, D.T., Singh, T., Kosciolek, T., and Tetchner, S. (2014). MetaPSI-COV: combining coevolution methods for accurate prediction of contacts and long range hydrogen bonding in proteins. Bioinformatics 31, 999–1006.
- 147. Xu, J. (2019). Distance-based protein folding powered by deep learning. Proc. Natl. Acad. Sci. U S A 116, 16856–16865.
- 148. Jones, D.T., and Kandathil, S.M. (2018). High precision in protein contact prediction using fully convolutional neural networks and minimal sequence features. Bioinformatics 34, 3308–3315.
- 149. Hanson, J., Paliwal, K., Litfin, T., Yang, Y., and Zhou, Y. (2018). Accurate prediction of protein contact maps by coupling residual two-dimensional bidirectional long short-term memory with convolutional neural networks. Bioinformatics 34, 4039–4045.
- 150. Kandathil, S.M., Greener, J.G., and Jones, D.T. (2019). Prediction of interresidue contacts with DeepMetaPSICOV in CASP13. Proteins 87, 1092–1099.
- 151. Hou, J., Wu, T., Cao, R., and Cheng, J. (2019). Protein tertiary structure modeling driven by deep learning and contact distance prediction in CASP13. Proteins 87, 1165–1178.
- 152. Zheng, W., Li, Y., Zhang, C., Pearce, R., Mortuza, S., and Zhang, Y. (2019). Deep-learning contact-map guided protein structure prediction in CASP13. Proteins 87, 1149–1164.
- 153. Wu, Q., Peng, Z., Anishchenko, I., Cong, Q., Baker, D., and Yang, J. (2020). Protein contact prediction using metagenome sequence data and residual neural networks. Bioinformatics 36, 41–48.
- 154. Marks, D.S., Colwell, L.J., Sheridan, R., Hopf, T.A., Pagnani, A., Zecchina, R., and Sander, C. (2011). Protein 3D structure computed from evolutionary sequence variation. PLoS One 6, e28766.
- 155. Ma, J., Wang, S., Wang, Z., and Xu, J. (2015). Protein contact prediction by integrating joint evolutionary coupling analysis and supervised learning. Bioinformatics *31*, 3506–3513.
- 156. Remmert, M., Biegert, A., Hauser, A., and Söding, J. (2012). HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment. Nat. Methods 9, 173–175.
- 157. Fariselli, P., Olmea, O., Valencia, A., and Casadio, R. (2001). Prediction of contact maps with neural networks and correlated mutations. Protein Eng. 14, 835–843.
- Horner, D.S., Pirovano, W., and Pesole, G. (2007). Correlated substitution analysis and the prediction of amino acid structural contacts. Brief. Bioinform. 9, 46–56.
- 159. Monastyrskyy, B., d'Andrea, D., Fidelis, K., Tramontano, A., and Kryshtafovych, A. (2014). Evaluation of residue–residue contact prediction in CASP10. Proteins 82, 138–153.
- 160. Xu, J., and Wang, S. (2019). Analysis of distance-based protein structure prediction by deep learning in CASP13. Proteins *87*, 1069–1081.
- 161. Moult, J., Fidelis, K., Kryshtafovych, A., Schwede, T., and Tramontano, A. (2018). Critical assessment of methods of protein structure prediction (CASP)—Round XII. Proteins 86, 7–15.
- Wang, S., Li, W., Liu, S., and Xu, J. (2016). RaptorX-Property: a web server for protein structure property prediction. Nucleic Acids Res. 44, W430–W435.
- 163. Gao, Y., Wang, S., Deng, M., and Xu, J. (2018). RaptorX-Angle: real-value prediction of protein backbone dihedral angles through a hybrid method of clustering and deep learning. BMC Bioinformatics 19, 100.



- 164. AlQuraishi, M. (2019). AlphaFold at CASP13. Bioinformatics 35, 4862–4865.
- Zemla, A., Venclovas, Č., Moult, J., and Fidelis, K. (1999). Processing and analysis of CASP3 protein structure predictions. Proteins 37, 22–29.
- 166. Kingma, D.P., Mohamed, S., Rezende, D.J., and Welling, M. (2014). Semi-supervised learning with deep generative models. Adv. Neural Inf. Process. Syst. 3581–3589.
- 167. Desmet, J., De Maeyer, M., Hazes, B., and Lasters, I. (1992). The deadend elimination theorem and its use in protein side-chain positioning. Nature 356, 539–542.
- Krivov, G.G., Shapovalov, M.V., and Dunbrack, R.L. (2009). Improved prediction of protein side-chain conformations with SCWRL4. Proteins 77, 778–795.
- 169. Liu, K., Sun, X., Ma, J., Zhou, Z., Dong, Q., Peng, S., Wu, J., Tan, S., Blobel, G., and Fan, J. (2017). Prediction of amino acid side chain conformation using a deep neural network. arXiv 1707, 08381.
- Du, Y., Meier, J., Ma, J., Fergus, R., and Rives, A. (2020). Energy-based models for atomic-resolution protein conformations. arXiv 2004, 13167.
- 171. LeCun, Y., Chopra, S., Hadsell, R., Ranzato, M., and Huang, F. (2006). A Tutorial on Energy-Based Learning (Predicting Structured Data), p. 1.
- 172. Zeng, H., Wang, S., Zhou, T., Zhao, F., Li, X., Wu, Q., and Xu, J. (2018). ComplexContact: a web server for inter-protein contact prediction using deep learning. Nucleic Acids Res. 46, W432–W437.
- 173. Wang, S., Li, Z., Yu, Y., and Xu, J. (2017). Folding membrane proteins by deep transfer learning. Cell Syst. 5, 202–211.e3.
- 174. Tsirigos, K.D., Peters, C., Shu, N., Käll, L., and Elofsson, A. (2015). The TOPCONS web server for consensus prediction of membrane protein topology and signal peptides. Nucleic Acids Res. 43, W401–W407.
- 175. Alford, R.F., and Gray, J.J. (2020). Big data from sparse data: diverse scientific benchmarks reveal optimization imperatives for implicit membrane energy functions. Biophys. J. 118, 361a.
- 176. Stein, A., and Kortemme, T. (2013). Improvements to robotics-inspired conformational sampling in Rosetta. PLoS One 8, e63090.
- 177. Ruffolo, J.A., Guerra, C., Mahajan, S.P., Sulam, J., and Gray, J.J. (2020). Geometric potentials from deep learning improve prediction of CDR H3 loop structures. Bioinformatics 36, i268–i275.
- 178. Nguyen, S.P., Li, Z., Xu, D., and Shang, Y. (2017). New deep learning methods for protein loop modeling. IEEE/ACM Trans. Comput. Biol. Bioinform. 16, 596–606.
- Li, Z.; Nguyen, S.P.; Xu, D.; Shang, Y. Protein loop modeling using deep generative adversarial network. Proceedings—International Conference on Tools with Artificial Intelligence, ICTAI. 2018; pp 1085–1091.
- 180. Porebski, B.T., and Buckle, A.M. (2016). Consensus protein design. Protein Eng. Des. Select. 29, 245–251.
- 181. Killoran, N., Lee, L.J., Delong, A., Duvenaud, D., and Frey, B.J. (2017). Generating and designing DNA with deep generative models. arXiv 1712, 06148.
- 182. Gupta, A., and Zou, J. (2018). Feedback GAN FBGAN for DNA: a novel feedback-loop architecture for optimizing protein functions. arXiv 1804, 01694.
- 183. Brookes, D.H., Park, H., and Listgarten, J. (2019). Conditioning by adaptive sampling for robust design. arXiv 1901, 10060.
- 184. Yu, C.-H., Qin, Z., Martin-Martinez, F.J., and Buehler, M.J. (2019). A selfconsistent sonification method to translate amino acid sequences into musical compositions and application in protein design using artificial intelligence. ACS Nano 13, 7471–7482.
- 185. Costello, Z., and Martin, H.G. (2019). How to hallucinate functional proteins. arXiv 1903, 00458.
- 186. Chhibbar, P., and Joshi, A. (2019). Generating protein sequences from antibiotic resistance genes data using generative adversarial networks. arXiv 1904, 13240.



- 187. Riesselman, A.J., Shin, J.-E., Kollasch, A.W., McMahon, C., Simon, E., Sander, C., Manglik, A., Kruse, A.C., and Marks, D.S. (2019). Accelerating protein design using autoregressive generative models. bioRxiv, 757252.
- Davidsen, K., Olson, B.J., DeWitt, W.S., III, Feng, J., Harkins, E., Bradley, P., and Matsen IV, F.A. (2019). Deep generative models for T cell receptor protein sequences. eLife 8, https://doi.org/10.7554/eLife.46935.
- 189. Han, X., Zhang, L., Zhou, K., and Wang, X. (2019). ProGAN: protein solubility generative adversarial nets for data augmentation in DNN framework. Comput. Chem. Eng. 131, 106533.
- Repecka, D., Jauniskis, V., Karpus, L., Rembeza, E., Zrimec, J., Poviloniene, S., et al. (2019). Expanding functional protein sequence space using generative adversarial networks. bioRxiv, 789719, https://doi.org/10. 1101/789719. https://www.biorxiv.org/content/10.1101/789719v2.
- Sabban, S., and Markovsky, M. (2020). RamaNet: computational de novo helical protein backbone design using a long short-term memory generative neural network. F1000Research 9, 298.
- 192. Eguchi, R.R., Anand, N., Choe, C.A., and Huang, P.-S. (2020). Ig-VAE: generative modeling of immunoglobulin proteins by direct 3D coordinate generation. bioRxiv, 242347. https://www.biorxiv.org/content/10.1101/ 2020.08.07.242347v1.
- Anishchenko, I., Chidyausiku, T.M., Ovchinnikov, S., Pellock, S.J., Baker, D., and Harvard, J. (2020). De novo protein design by deep network hallucination. bioRxiv, 211482. https://www.biorxiv.org/content/10.1101/ 2020.07.22.211482v1.
- 194. Wang, J., Cao, H., Zhang, J.Z., and Qi, Y. (2018). Computational protein design with deep learning neural networks. Sci. Rep. 8, 6349.
- 195. Greener, J.G., Moffat, L., and Jones, D.T. (2018). Design of metalloproteins and novel protein folds using variational autoencoders. Sci. Rep. 8, 1–12.
- 196. Chen, S., Sun, Z., Lin, L., Liu, Z., Liu, X., Chong, Y., Lu, Y., Zhao, H., and Yang, Y. (2019). To improve protein sequence profile prediction through image captioning on pairwise residue distance map. J. Chem. Inf. Model. 60, 391–399.
- 197. Zhang, Y., Chen, Y., Wang, C., Lo, C.-C., Liu, X., Wu, W., and Zhang, J. (2019). ProDCoNN: protein design using a convolutional neural network. Proteins 88, 819–829.
- 198. Shroff, R., Cole, A.W., Morrow, B.R., Diaz, D.J., Donnell, I., Gollihar, J., Ellington, A.D., and Thyer, R. (2019). A structure-based deep learning framework for protein engineering. bioRxiv, 833905.
- Strokach, A., Becerra, D., Corbi-Verge, C., Perez-Riba, A., and Kim, P.M. (2019). Designing real novel proteins using deep graph neural networks. bioRxiv, 868935.
- 200. Karimi, M., Zhu, S., Cao, Y., and Shen, Y. (2019). De novo protein design for novel folds using guided conditional Wasserstein generative adversarial networks gcWGAN. bioRxiv, 769919.
- 201. Qi, Y., and Zhang, J.Z. (2020). DenseCPD: improving the accuracy of neural-network-based computational protein sequence design with DenseNet. J. Chem. Inf. Model. 60, 1245–1252.
- 202. Anand, N., Eguchi, R.R., Derry, A., Altman, R.B., and Huang, P. (2020). Protein sequence design with a learned potential. bioRxiv, 895466.
- Norn, C., Wicky, B.I., Juergens, D., Liu, S., Kim, D., Koepnick, B., et al. (2020). Protein sequence design by explicit energy landscape optimization. bioRxiv, 218917, https://doi.org/10.1101/2020.07.23.218917. https://www.biorxiv.org/content/10.1101/2020.07.23.218917v1.full.
- 204. Waghu, F.H., Gopi, L., Barai, R.S., Ramteke, P., Nizami, B., and Idicula-Thomas, S. (2014). CAMP: collection of sequences and structures of antimicrobial peptides. Nucleic Acids Res. 42, D1154–D1158.
- 205. Grisoni, F., Neuhaus, C.S., Gabernet, G., Müller, A.T., Hiss, J.A., and Schneider, G. (2018). Designing anticancer peptides by constructive machine learning. ChemMedChem 13, 1300–1302.
- 206. Yu, F., and Koltun, V. (2015). Multi-scale context aggregation by dilated convolutions. arXiv 1511, 07122.



- 207. Gupta, A., and Zou, J. (2019). Feedback GAN for DNA optimizes protein functions. Nat. Machine Intelligence 1, 105–111.
- Kuhlman, B., and Baker, D. (2000). Native protein sequences are close to optimal for their structures. Proc. Natl. Acad. Sci. U S A 97, 10383–10388.
- 209. Li, Z., Yang, Y., Faraggi, E., Zhan, J., and Zhou, Y. (2014). Direct prediction of profiles of sequences compatible with a protein structure by neural networks with fragment-based local and energy-based nonlocal profiles. Proteins 82, 2565–2573.
- Karimi, M., Zhu, S., Cao, Y., and Shen, Y. (2020). De novo protein design for novel folds using guided conditional Wasserstein generative adversarial networks. J. Chem. Inf. Model. <u>https://doi.org/10.1021/acs.jcim.</u> 0c00593.
- 211. Hou, J., Adhikari, B., and Cheng, J. (2017). DeepSF: deep convolutional neural network for mapping protein sequences to folds. Bioinformatics 34, 1295–1303.
- 212. Jelinek, F., Mercer, R.L., Bahl, L.R., and Baker, J.K. (1977). Perplexity a measure of the difficulty of speech recognition tasks. J. Acoust. Soc. Am. 62, S63.
- Strokach, A., Becerra, D., Corbi-Verge, C., Perez-Riba, A., and Kim, P. (2019). Fast and flexible design of novel proteins using graph neural networks. bioRxiv, 868935.
- 214. Ramachandran, G.N. (1963). Stereochemistry of polypeptide chain configurations. J. Mol. Biol. 7, 95–99.
- (2015). https://research.googleblog.com/2015/06/Ωinceptionism-goingdeeper-into-neural.html.
- 216. Sutton, R.S., and Barto, A.G. (2018). Reinforcement Learning: An Introduction (MIT press).
- Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; Fei-Fei, L. Imagenet: A large-scale hierarchical image database. 2009 IEEE Conference on Computer Vision and Pattern Recognition 2009, 248–255.
- Mayr, A., Klambauer, G., Unterthiner, T., and Hochreiter, S. (2016). Deep-Tox: toxicity prediction using deep learning. Front. Environ. Sci. 3, 80.
- Brown, N., Fiscato, M., Segler, M.H., and Vaucher, A.C. (2019). Guaca-Mol: benchmarking models for de novo molecular design. J. Chem. Inf. Model. 59, 1096–1108.
- 220. Lutter, M., Ritter, C., and Peters, J. (2019). Deep Lagrangian networks: using physics as model prior for deep learning. arXiv 1907, 04490.
- 221. Greydanus, S., Dzamba, M., and Yosinski, J. (2019). Hamiltonian neural networks. Adv. Neural Inf. Process. Syst. 15379–15389.
- 222. Raissi, M., Perdikaris, P., and Karniadakis, G.E. (2019). Physics-informed neural networks: a deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. J. Comput. Phys. 378, 686–707.
- 223. Zepeda-Núñez, L., Chen, Y., Zhang, J., Jia, W., Zhang, L., and Lin, L. (2019). Deep Density: circumventing the Kohn-Sham equations via symmetry preserving neural networks. arXiv 1912, 00775.
- 224. Han, J., Li, Y., Lin, L., Lu, J., Zhang, J., and Zhang, L. (2019). Universal approximation of symmetric and anti-symmetric functions. arXiv 1912, 01765.
- 225. Shapovalov, M.V., and Dunbrack, R.L., Jr. (2011). A smoothed backbone-dependent rotamer library for proteins derived from adaptive kernel density estimates and regressions. Structure *19*, 844–858.
- Hintze, B.J., Lewis, S.M., Richardson, J.S., and Richardson, D.C. (2016). Molprobity's ultimate rotamer-library distributions for model validation. Proteins 84, 1177–1189.
- 227. Jensen, K.F., Coley, C.W., and Eyke, N.S. (2019). Autonomous discovery in the chemical sciences part I: progress. Angew. Chem. Int. Ed. 59, 2–38.
- Coley, C.W., Eyke, N.S., and Jensen, K.F. (2019). Autonomous discovery in the chemical sciences part II: outlook. Angew. Chem. Int. Ed. 59, 2–25.
- 229. Coley, C.W., Thomas, D.A., Lummiss, J.A., Jaworski, J.N., Breen, C.P., Schultz, V., Hart, T., Fishman, J.S., Rogers, L., Gao, H., et al. (2019). A



robotic platform for flow synthesis of organic compounds informed by Al planning. Science 365, eaax1566.

- Barrett, R.; White, A.D. Iterative peptide modeling with active learning and meta-learning. arXiv preprint 2019, 1911.09103.
- 231. You, J., Liu, B., Ying, R., Pande, V., and Leskovec, J. (2018). Graph convolutional policy network for goal-directed molecular graph generation. Adv. Neural Inf. Process. Syst. 6410–6421.
- 232. Zhou, Z., Kearnes, S., Li, L., Zare, R.N., and Riley, P. (2019). Optimization of molecules via deep reinforcement learning. Sci. Rep. 9, 1–10.
- 233. Mirhoseini, A., Goldie, A., Yazgan, M., Jiang, J., Songhori, E., Wang, S., Lee, Y.-J., Johnson, E., Pathak, O., Bae, S., et al. (2004). Chip placement with deep reinforcement learning. arXiv 2020, 10746.
- 234. Cooper, S., Khatib, F., Treuille, A., Barbero, J., Lee, J., Beenen, M., Leaver-Fay, A., Baker, D., Popović, Z., and Foldit, P. (2010). Predicting protein structures with a multiplayer online game. Nature 466, 756–760.
- 235. Koepnick, B., Flatten, J., Husain, T., Ford, A., Silva, D.-A., Bick, M.J., Bauer, A., Liu, G., Ishida, Y., Boykov, A., et al. (2019). De novo protein design by citizen scientists. Nature 570, 390–394.
- 236. Czibula, G., Bocicor, M.-I., and Czibula, I.-G. (2011). A reinforcement learning model for solving the folding problem. Int. J. Comput. Technol. Appl. 2, 171–182.
- 237. Jafari, R., and Javidi, M.M. (2020). Solving the protein folding problem in hydrophobic-polar model using deep reinforcement learning. SN Appl. Sci. 2, 259.
- 238. Gao, W. (2020). Development of a Protein Folding Environment for Reinforcement Learning, M.Sc. thesis (Johns Hopkins University).
- Angermueller, C., Dohan, D., Belanger, D., Deshpande, R., Murphy, K., and Colwell, L. (2020). Model-Based Reinforcement Learning for Biological Sequence Design (ICLR 2020 Conference). https://openreview.net/ forum?id=HklxbaBKvr.
- 240. Zeiler, M.D., and Fergus, R. (2014). Visualizing and understanding convolutional networks. Eur. Conf. Comput. Vis. 818–833.
- 241. Smilkov, D., Thorat, N., Kim, B., Viégas, F., and Wattenberg, M. (2017). SmoothGrad: removing noise by adding noise. arXiv 1706, 03825.

- Sundararajan, M.; Taly, A.; Yan, Q. Axiomatic attribution for deep networks. Proceedings of the 34th International Conference on Machine Learning2017, 70, 3319–3328.
- 243. Adebayo, J., Gilmer, J., Muelly, M., Goodfellow, I., Hardt, M., and Kim, B. (2018). Sanity checks for saliency maps. Adv. Neural Inf. Process. Syst. 9505–9515.
- 244. Shrikumar, A., Greenside, P., and Kundaje, A. (1704). Learning important features through propagating activation differences. arXiv 2017, 02685.
- Lundberg, S.M.; Lee, S.-I. A Unified Approach to Interpreting Model Predictions. Proceedings of the 31st International Conference on Neural Information Processing Systems 2017, 4768–4777.
- 246. Hannon, G.J. (2002). RNA interference. Nature 418, 244-251.
- 247. Zhang, P., Woen, S., Wang, T., Liau, B., Zhao, S., Chen, C., Yang, Y., Song, Z., Wormald, M.R., Yu, C., et al. (2016). Challenges of glycosylation analysis and control: an integrated approach to producing optimal and consistent therapeutic drugs. Drug Discov. Today *21*, 740–765.
- Sanchez-Lengeling, B., and Aspuru-Guzik, A. (2018). Inverse molecular design using machine learning: generative models for matter engineering. Science 361, 360–365.
- 249. Coley, C.W., Jin, W., Rogers, L., Jamison, T.F., Jaakkola, T.S., Green, W.H., Barzilay, R., and Jensen, K.F. (2019). A graph-convolutional neural network model for the prediction of chemical reactivity. Chem. Sci. 10, 370–377.
- 250. Yang, K., Swanson, K., Jin, W., Coley, C., Eiden, P., Gao, H., Guzman-Perez, A., Hopper, T., Kelley, B., Mathea, M., et al. (2019). Analyzing learned molecular representations for property prediction. J. Chem. Inf. Model. 59, 3370–3388.
- Gao, W., and Coley, C.W. (2020). The synthesizability of molecules proposed by generative models. J. Chem. Inf. Model. <u>https://doi.org/10.1021/acs.jcim.0c00174</u>.
- 252. Langan, R.A., Boyken, S.E., Ng, A.H., Samson, J.A., Dods, G., Westbrook, A.M., Nguyen, T.H., Lajoie, M.J., Chen, Z., Berger, S., et al. (2019). De novo design of bioactive protein switches. Nature 572, 205–210.