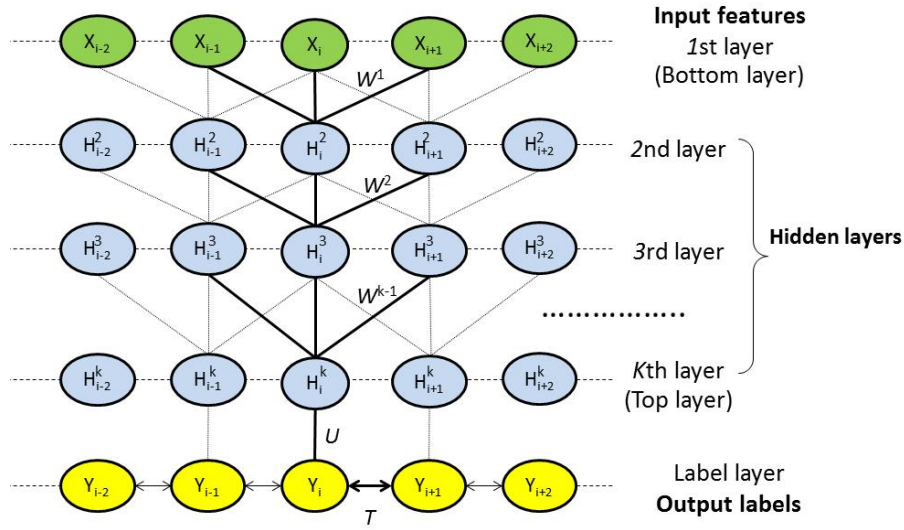


## Supplemental Document

### DeepCNF model

As shown in Supplemental Figure 1, DeepCNF consists of two modules: (a) the Conditional Random Fields (CRF) [1] module consisting of the top layer and the label layer, and (b) the deep convolutional neural network (DCNN) [2] module covering the input to the top layer. When only one hidden layer is used, this DeepCNF becomes Conditional Neural Fields (CNF), a probabilistic graphical model described in [3].



**Supplemental Figure 1.** The architecture of DeepCNF, where  $i$  is the residue index,  $X_i$  the associated input features,  $H^k$  represents the  $k$ th hidden layer, and  $Y$  is the output label. All the layers from the 1st to the top layer form a deep convolutional neural network (DCNN). The top layer and the label layer form a conditional random field (CRF).  $W^k \{k = 1, 2, \dots, K\}$ ,  $U$  and  $T$  are the model parameters where  $T$  is used to model correlation among adjacent residues.

**Conditional Random Field (CRF).** Given a protein sequence of length  $L$ , let  $\mathbf{Y} = (Y_1, \dots, Y_L)$  denote its SS where  $Y_i$  is the SS type at residue  $i$ . Let  $\mathbf{X} = (X_1, \dots, X_L)$  denote the input feature where  $X_i$  is a column vector representing the input feature for residue  $i$ . Using DeepCNF, we calculate the conditional probability of  $\mathbf{Y}$  on the input  $\mathbf{X}$  as follows,

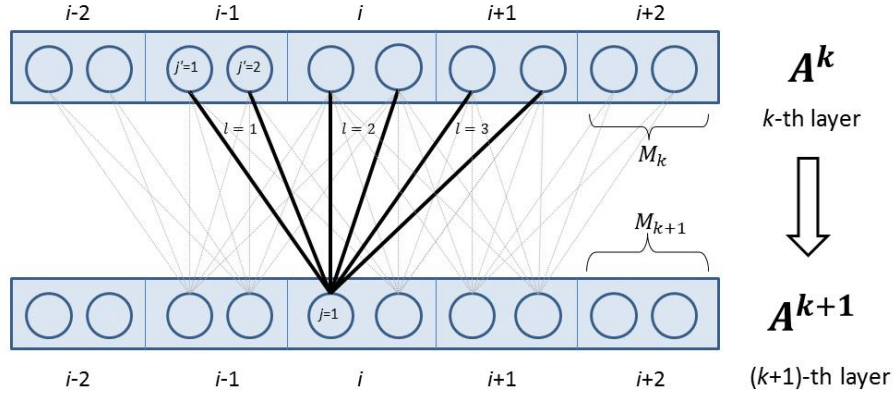
$$P(\mathbf{Y}|\mathbf{X}) = \exp(\sum_{i=1}^L [\Psi(\mathbf{Y}, \mathbf{X}, i) + \Phi(\mathbf{Y}, \mathbf{X}, i)]) / Z(\mathbf{X}) \quad (1)$$

Where  $\Psi(\mathbf{Y}, \mathbf{X}, i)$  is the potential function quantifying correlation among adjacent SS types at around position  $i$ ,  $\Phi(\mathbf{Y}, \mathbf{X}, i)$  is the potential function modeling relationship between  $Y_i$  and input features for position  $i$ , and  $Z(\mathbf{X})$  is the partition function. Formally,  $\Psi()$  and  $\Phi()$  are defined as follows,

$$\Psi(\mathbf{Y}, \mathbf{X}, i) = \sum_{a,b} T_{a,b} \delta(Y_i = a) \delta(Y_{i+1} = b) \quad (2)$$

$$\Phi(\mathbf{Y}, \mathbf{X}, i) = \sum_a \sum_m U_{a,m} H_m(\mathbf{X}, i, W) \delta(Y_i = a) \quad (3)$$

Where  $a$  and  $b$  represent secondary structure states,  $\delta()$  is an indicator function,  $H_m(\mathbf{X}, i, W)$  is a neural network function for the  $m$ -th neuron at position  $i$  of the top layer, and  $W$ ,  $U$ , and  $T$  are the model parameters to be trained. Specifically,  $W$  is the parameter for the neural network,  $U$  is the parameter connecting the top layer to the label layer, and  $T$  is for label correlation. Below see the details of the deep convolutional neural network for  $H_m(\mathbf{X}, i, W)$ .

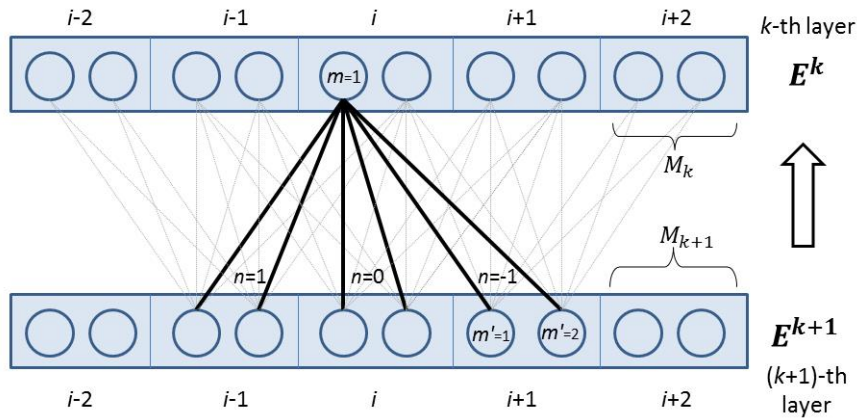


**Supplemental Figure 2.** The feed-forward connection between two adjacent layers in the deep convolutional neural network.

**Deep convolutional neural network (DCNN).** Supplemental Figure 2 shows two adjacent layers. Let  $M_k$  be the number of neurons for a single position at the  $k$ -th layer. Let  $X_i(m)$  be the  $m$ -th feature at the input layer for residue  $i$  and  $A_i^k(m)$  denote the output value of the  $m$ -th neuron of position  $i$  at layer  $k$ . When  $k = 1$ ,  $A^k$  is actually the input feature  $X$ . Otherwise,  $A^k$  is a matrix of dimension  $L \times M_k$ . Let  $2N_k + 1$  be the window size at the  $k$ -th layer. Mathematically,  $A_i^k(m)$  is defined as follows.

$$\begin{aligned}
 A_i^k(m) &= X_i(m), & \text{if } k = 1 \\
 A_i^{k+1}(m) &= h\left(\sum_{n=-N_k}^{N_k} \sum_{m'=1}^{M_k} [A_{i+n}^k(m') * W_n^k(m, m')]\right), & \text{if } k < K \\
 H_m(X, i, W) &= A_i^k(m) & \text{if } k = K
 \end{aligned} \tag{4}$$

Meanwhile,  $h()$  is the activation function, either the sigmoid (i.e.,  $1/(1 + \exp(-x))$ ) or the tanh (i.e.,  $(1 - \exp(-2x))/(1 + \exp(-2x))$ ) function.  $W_n^k$  where  $(-N_k \leq n \leq N_k)$  is a 2D weight matrix for the connections between the neurons of position  $i + n$  at layer  $k$  and the neurons of position  $i$  at layer  $k + 1$ .  $W_n^k(m, m')$  is shared by all the positions in the same layer, so it is position-independent. Here  $m'$  and  $m$  index two neurons at the  $k$ -th and  $(k + 1)$ -th layers, respectively.



**Supplemental Figure 3.** Illustration of calculating the gradient of deep convolutional neural network from layer  $k + 1$  to layer  $k$ .

## Training method

Similar to CRF [1], we train the model parameters by maximum-likelihood. The log-likelihood is as follows.

$$\log P(\mathbf{Y}|\mathbf{X}) = \sum_{i=1}^L [\Psi(\mathbf{Y}, \mathbf{X}, i) + \Phi(\mathbf{Y}, \mathbf{X}, i)] - \log Z(\mathbf{X}) \quad (5)$$

To train the model parameters, we need to calculate the gradient with respect to each parameter. We calculate the gradient first for CRF and then for DCNN. The gradient of the log-likelihood with respect to the parameters  $T$  and  $U$  is given by,

$$\nabla_{T_{a,b}} = [\sum_{i=1}^L \delta(Y_i = a) \delta(Y_{i+1} = b)] - EXP_{P(\tilde{\mathbf{Y}}|\mathbf{X}, \mathbf{W}, \mathbf{U}, T)} [\sum_{i=1}^L \delta(\tilde{Y}_i = a) \delta(\tilde{Y}_{i+1} = b)] \quad (6)$$

$$\nabla_{U_{a,m}} = [\sum_{i=1}^L \delta(Y_i = a) H_m(\mathbf{X}, i, \mathbf{W})] - EXP_{P(\tilde{\mathbf{Y}}|\mathbf{X}, \mathbf{W}, \mathbf{U}, T)} [\sum_{i=1}^L \delta(\tilde{Y}_i = a) H_m(\mathbf{X}, i, \mathbf{W})] \quad (7)$$

Where  $EXP$  is the expectation function and can be calculated efficiently using the forward-backward algorithm [3].

As shown in Supplemental Figure 3, we can calculate the neuron error values at the  $k$ -th layer in a back-propagation mode as follows.

$$\begin{aligned} E_i^k(m) &= g(A_i^k(m)) * \sum_a [E_i(a) * U_{a,m}] \quad , \quad \text{if } k = K \\ \text{where } E_i(a) &= [\delta(Y_i = a)] - EXP_{P(\tilde{\mathbf{Y}}|\mathbf{X}, \mathbf{W}, \mathbf{U}, T)} [\delta(\tilde{Y}_i = a)] \\ E_i^k(m) &= g(A_i^k(m)) * \sum_{n=-N_k}^{N_k} \sum_{m'=1}^{M_{k+1}} [E_{i-n}^{k+1}(m') * W_n^k(m', m)], \quad \text{if } k < K \end{aligned} \quad (8)$$

Where  $g()$  is the derivative of the activation function; it is  $g(x) = (1-x)x$  and  $g(x) = 1-x^2$  for the sigmoid and tanh function, respectively.  $\mathbf{E}^k$  is the neuron error value matrix at the  $k$ -th layer, with dimension  $L \times M_k$ .  $E_i(a)$  is the error of the log-likelihood function with respect to the label at the  $i$ -th position and can be calculated by the forward-backward algorithm. Finally, the gradient of the parameter  $W$  at the  $k$ -th layer is:

$$\nabla_{W_n^k(m, m')} = \sum_{i=1}^L [E_i^{k+1}(m) * A_{i+n}^k(m')] \quad (9)$$

## L2 regularization and L-BFGS

To reduce over-fitting, the log-likelihood objective function is penalized with a  $L_2$ -norm of the model parameters. Thus, our final objective function is as follows.

$$\max_{\theta} \log P_{\theta}(\mathbf{Y}|\mathbf{X}) - \lambda \|\theta\|^2 \quad (10)$$

Where  $\theta$  is the set of model parameters and  $\lambda$  is the regularization factor used to avoid overfitting. Although DeepCNF has a large number of model parameters, by setting the regularization factor large enough, we can make the  $L_2$ -norm of the model parameters small and thus, restrict the search space of the model parameter and avoid overfitting. However, a very large regularization factor (e.g., infinity) may restrict the model parameter into too small a search space and the resultant model may not learn enough from the training data (i.e., under-fitting). We will determine the regularization factor by cross-validation.

Since the log-likelihood function is not convex, usually we can only solve the objective function to a local instead of global optimum. Although a typical way to train a deep network is to do it layer-by-layer, in our implementation we train all the model parameters simultaneously. We use the L-BFGS [4] to search for the optimal model parameters, which has also been successfully used to train CRF and CNF. In addition to learning the model parameters simultaneously, we can also train them

layer-by-layer in a supervised mode. Starting from the first layer (i.e., input feature), we train the model parameter  $W1$  by removing the third to the  $K$ -th layers but keeping the label layer. After  $W1$  is trained, we generate the neuron output values for the second layer and use them as input to train the model parameter  $W2$  by removing the fourth to the  $K$ -th layers but keeping the label layer. We repeat this procedure until all the parameters are trained, and finally we fine-tune these parameters by simultaneous training.

## References

1. Lafferty J, McCallum A, Pereira FC. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. 2001.
2. Lee H, Grosse R, Ranganath R, Ng AY. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. Proceedings of the 26th Annual International Conference on Machine Learning; 2009: ACM; 2009. p. 609-16.
3. Peng J, Bo L, Xu J. Conditional neural fields. Advances in neural information processing systems; 2009; 2009. p. 1419-27.
4. Liu DC, Nocedal J. On the limited memory BFGS method for large scale optimization. *Mathematical programming* 1989; **45**: 503-28.