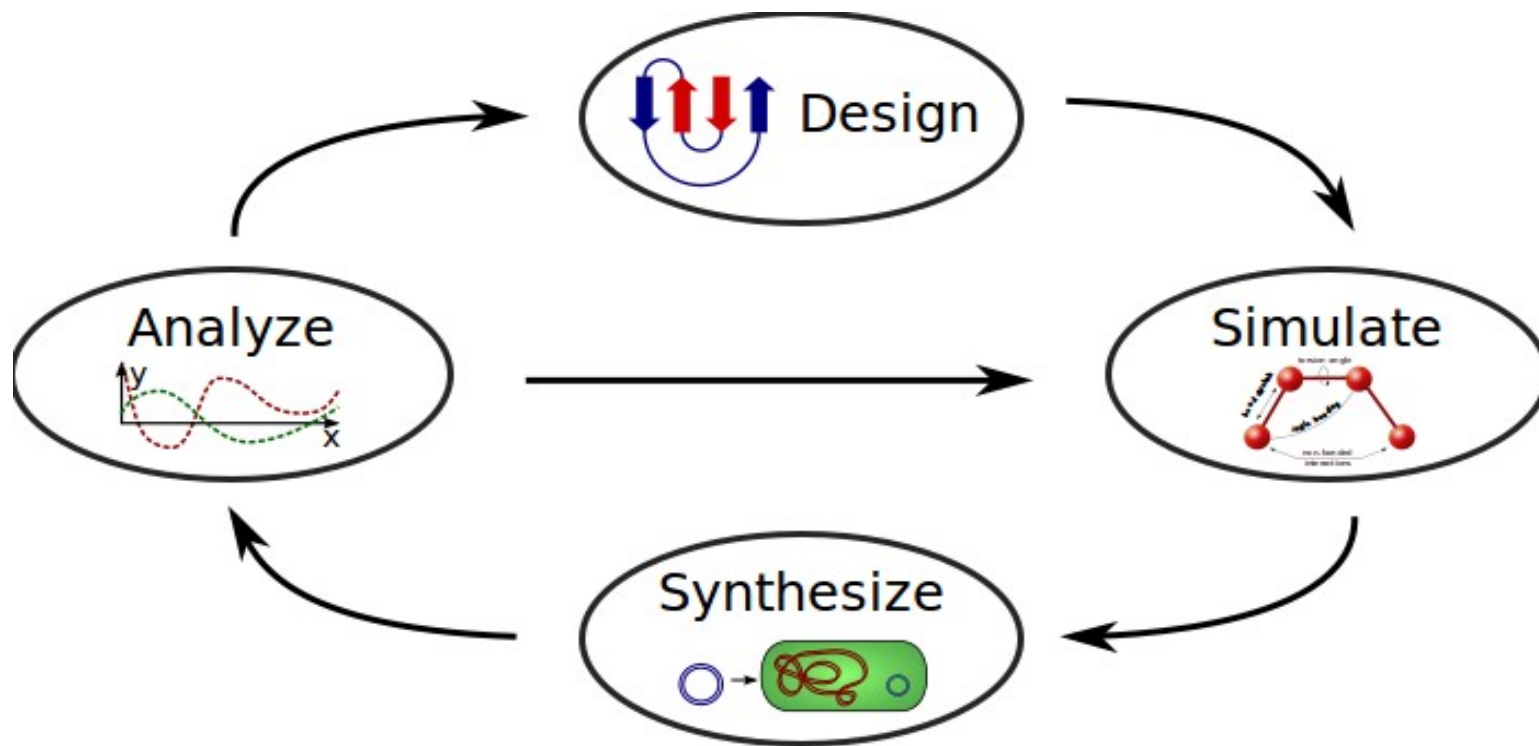


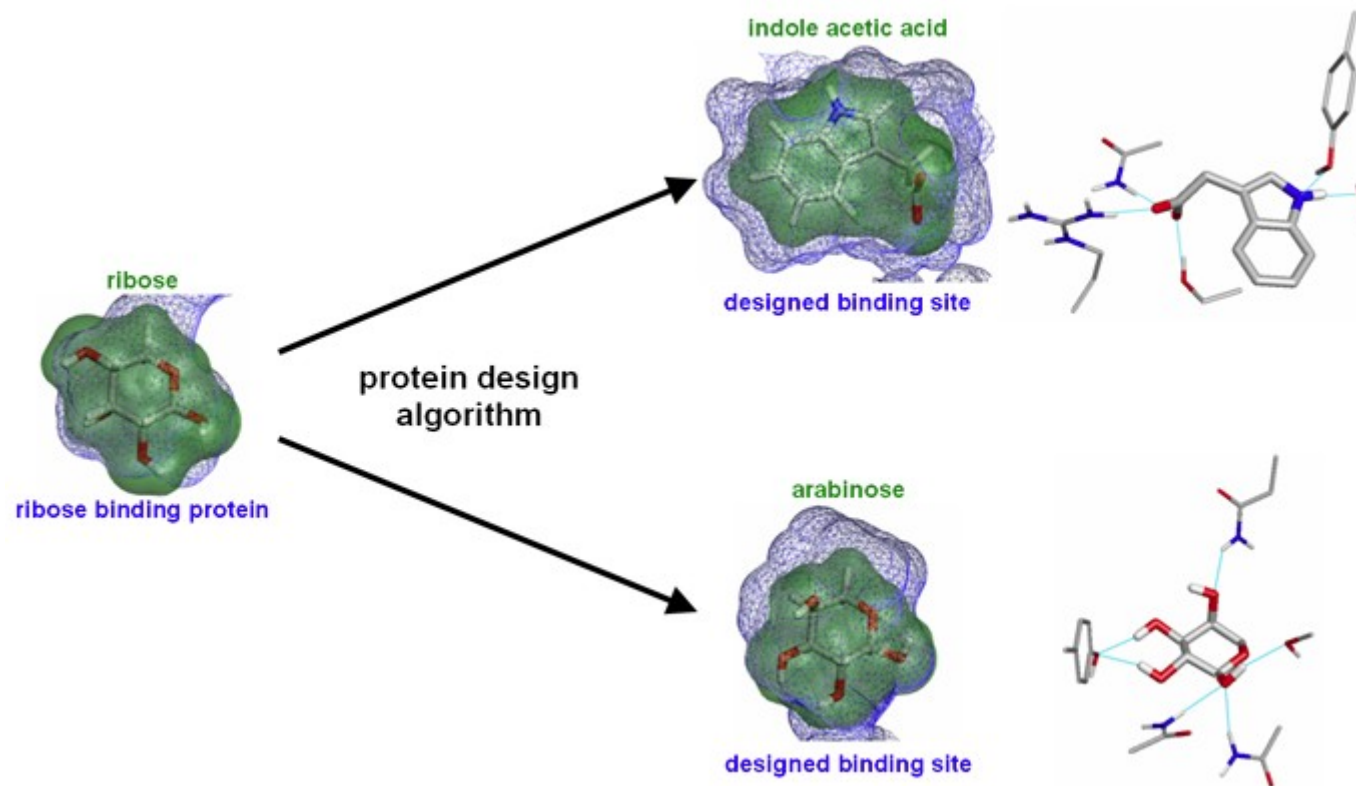
# Pareto-Optimal Multistate Protein Design Using Tree Graphical Models

- *Subhodeep Moitra*

# Protein Engineering Cycle



# Protein Design



Theory

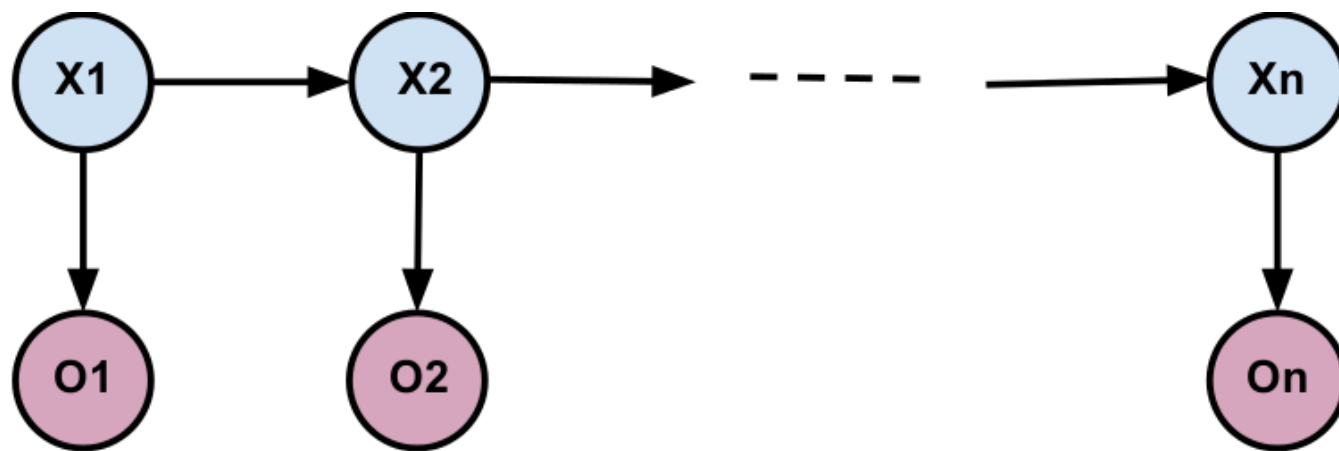


Figure 1: A canonical HMM

The graphical representation of the HMM factorizes according to the following probability distribution :

$$P(X_1, \dots, X_n | \mathbf{O}_1, \dots, \mathbf{O}_n) = P(X_1) \prod_{i=1}^n P(\mathbf{O}_i | X_i) \prod_{i=1}^{n-1} P(X_{i+1} | X_i) \quad (1)$$

*Calculation of the Optimal Probability value:*

$$V(i, k) = \begin{cases} P(\mathbf{O}_1|X_1=k)P(X_1=k) & \text{if } i = 1 \\ P(\mathbf{O}_i|X_i=k) \max_j \{P(X_i=k|X_{i-1}=j)V(i-1, j)\} & \text{if } 2 \leq i \leq n \end{cases}$$

*Retrieval of the optimal sequence corresponding to the optimal probability value:*

$$X_n^* = \arg \max_k V(n, k)$$

$$X_i^* = Ptr(X_{i+1}^*, i+1) \quad \text{if } 1 \leq i \leq n-1$$

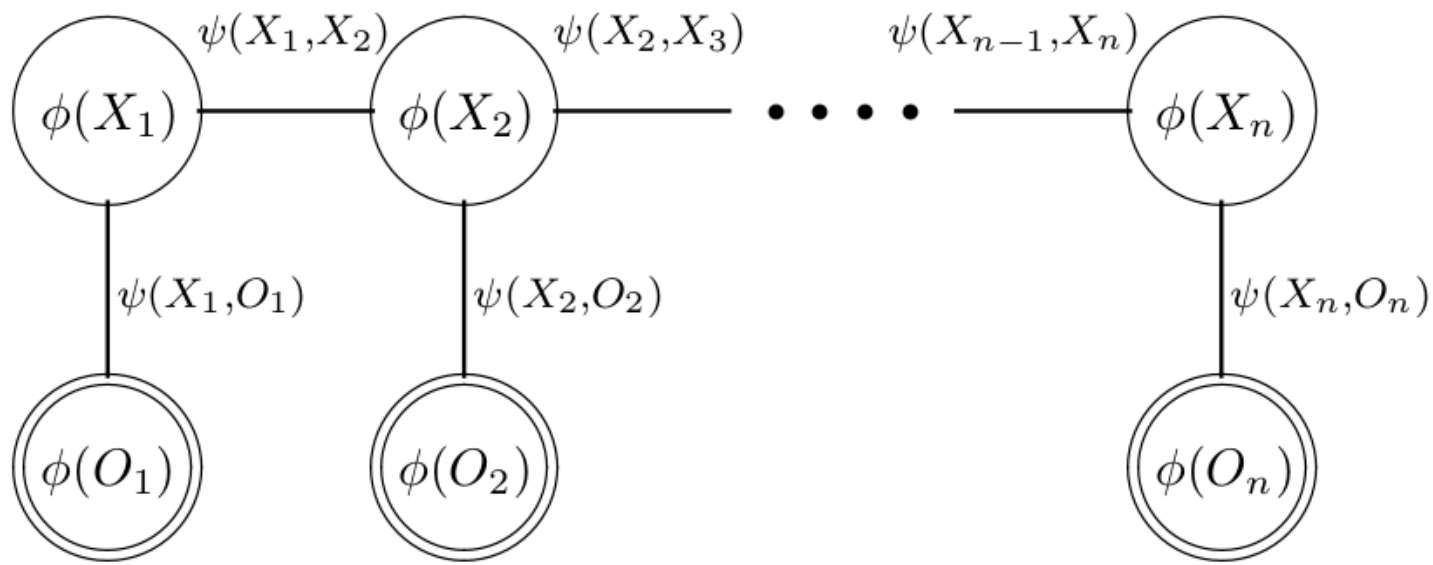


Figure 2: A chain structured Markov Random Field. Probability of a sequence assignment is defined as a function of the node and edge potentials

The probability of the a sequence assignment is defined as follows :

$$P(X_1, X_2, \dots, X_n) = \frac{1}{Z} \prod_{i=1}^n \phi(X_i) \phi(\mathbf{O}_i) \psi(\mathbf{O}_i, X_i) \prod_{i=1}^{n-1} \psi(X_i, X_{i+1})$$

where  $Z = \sum_{X_1, X_2, \dots, X_n} \left[ \prod_{i=1}^n \phi(X_i) \phi(\mathbf{O}_i) \psi(\mathbf{O}_i, X_i) \prod_{i=1}^{n-1} \psi(X_i, X_{i+1}) \right]$

The main steps in the max-product algorithm are :

*Calculation of the Optimal Probability value:*

$$V(i, k) = \begin{cases} \phi(\mathbf{O}_1)\psi(\mathbf{O}_1, X_1=k)\phi(X_1=k) & \text{if } i = 1 \\ \phi(\mathbf{O}_i)\phi(X_i=k)\psi(\mathbf{O}_i, X_i=k) \max_j \{\psi(X_i=k, X_{i-1}=j)V(i-1, j)\} & \text{if } 2 \leq i \leq n \end{cases}$$



The main steps in the max-product algorithm are :

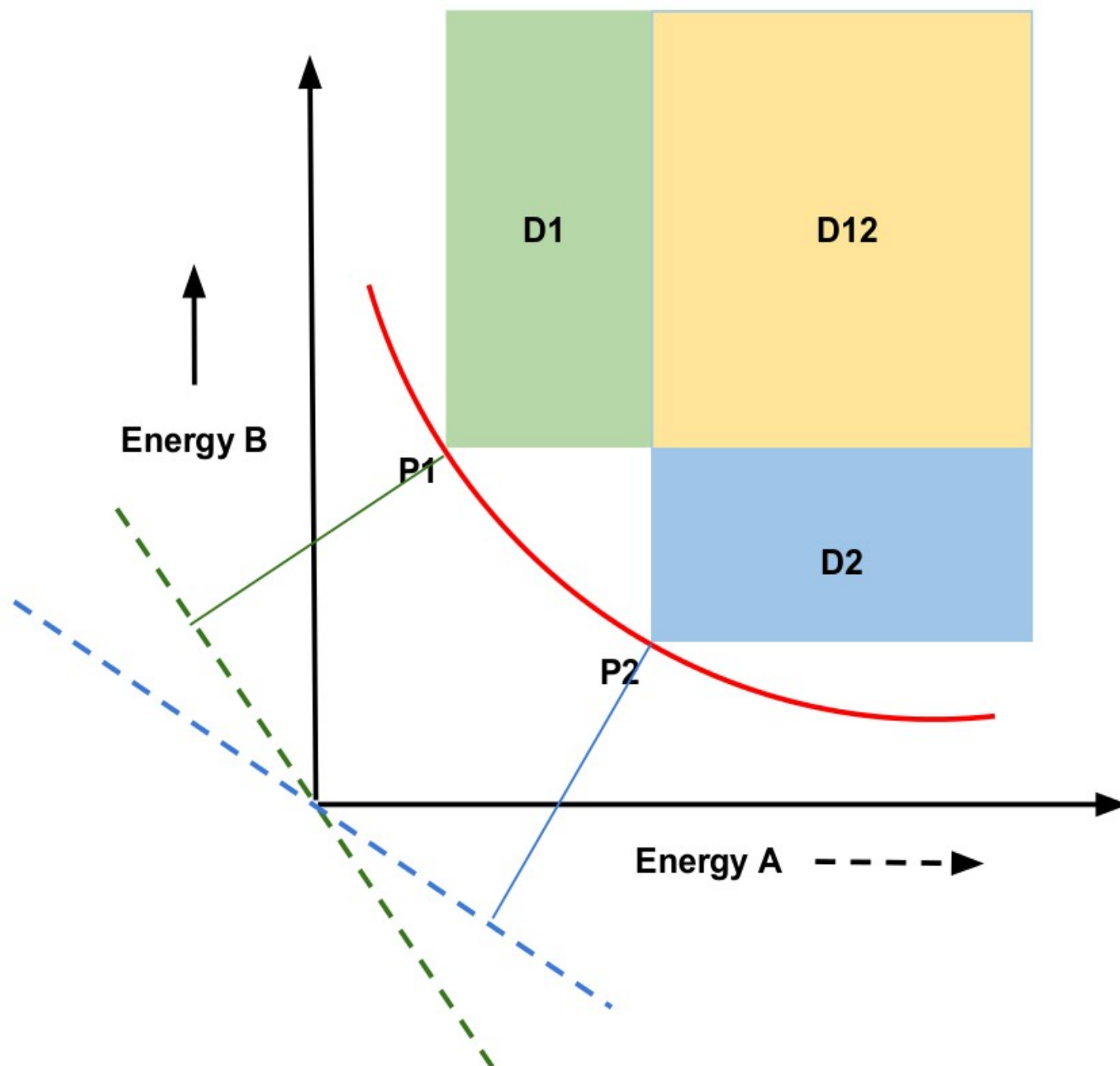
*Calculation of the Optimal Probability value:*

$$V(i, k) = \begin{cases} \phi(\mathbf{O}_1)\psi(\mathbf{O}_1, X_1=k)\phi(X_1=k) & \text{if } i = 1 \\ \phi(\mathbf{O}_i)\phi(X_i=k)\psi(\mathbf{O}_i, X_i=k) \max_j \{\psi(X_i=k, X_{i-1}=j)V(i-1, j)\} & \text{if } 2 \leq i \leq n \end{cases}$$

*Retrieval of the optimal sequence corresponding to the optimal probability value:*

$$X_n^* = \arg \max_k V(n, k)$$

$$X_i^* = Ptr(X_{i+1}^*, i+1) \quad \text{if } 1 \leq i \leq n-1$$



For multi-state protein design we can combine the energy functions as follows.

$$E_{AB}(X) = \theta_A E_A(X) + \theta_B E_B(X) \quad \text{where} \quad \theta_A, \theta_B \geq 0 \text{ and } \theta_A + \theta_B = 1$$

Here  $E_A(X)$  is any energy function mapping sequence to energy values for state A. The lower the energy the better the fit for sequence X to state A.  $E_A(X)$  can also be written as a CMRF.

$$E_A(X) = -\log \left[ \prod_{i=1}^n \phi(X_i) \phi(\mathbf{O}_i^{\mathbf{A}}) \psi(\mathbf{O}_i^{\mathbf{A}}, X_i) \prod_{i=1}^{n-1} \psi(X_i, X_{i+1}) \right]$$

$$X_A = \arg \min_X E_A(X)$$

$$\begin{aligned}
&= \arg \min_X - \log \left[ \prod_{i=1}^n \phi(X_i) \phi(\mathbf{O}_i^{\mathbf{A}}) \psi(\mathbf{O}_i^{\mathbf{A}}, X_i) \prod_{i=1}^{n-1} \psi(X_i, X_{i+1}) \right] \\
&= \arg \min_X - \left[ \sum_{i=1}^n \log \phi(X_i) + \sum_{i=1}^n \log \phi(\mathbf{O}_i^{\mathbf{A}}) + \sum_{i=1}^n \log \psi(\mathbf{O}_i^{\mathbf{A}}, X_i) + \sum_{i=1}^{n-1} \log \psi(X_i, X_{i+1}) \right] \\
&= \arg \min_X - \left[ \sum_{i=1}^n \log \phi(X_i) + \sum_{i=1}^n \log \psi(\mathbf{O}_i^{\mathbf{A}}, X_i) + \sum_{i=1}^{n-1} \log \psi(X_i, X_{i+1}) \right]
\end{aligned}$$

Similarly,

$$X_B = \arg \min_X E_B(X)$$

$$= \arg \min_X - \left[ \sum_{i=1}^n \log \phi(X_i) + \sum_{i=1}^n \log \psi(\mathbf{O}_i^{\mathbf{B}}, X_i) + \sum_{i=1}^{n-1} \log \psi(X_i, X_{i+1}) \right]$$

Taking the convex combination of the energy functions we get

$$\begin{aligned} X_{AB} &= \arg \min_X \theta_A E_A(X) + \theta_B E_B(X) \\ \text{s.t. } &\theta_A + \theta_B = 1 \\ &\theta_A, \theta_B \geq 0 \end{aligned}$$

Taking the convex combination of the energy functions we get

$$X_{AB} = \arg \min_X \theta_A E_A(X) + \theta_B E_B(X)$$

$$\text{s.t } \theta_A + \theta_B = 1$$

$$\theta_A, \theta_B \geq 0$$

$$X_{AB} = \arg \min_X -\theta_A \left[ \sum_{i=1}^n \log \phi(X_i) + \sum_{i=1}^n \log \psi(\mathbf{O}_i^{\mathbf{A}}, X_i) + \sum_{i=1}^{n-1} \log \psi(X_i, X_{i+1}) \right]$$

$$- \theta_B \left[ \sum_{i=1}^n \log \phi(X_i) + \sum_{i=1}^n \log \psi(\mathbf{O}_i^{\mathbf{B}}, X_i) + \sum_{i=1}^{n-1} \log \psi(X_i, X_{i+1}) \right]$$

Taking the convex combination of the energy functions we get

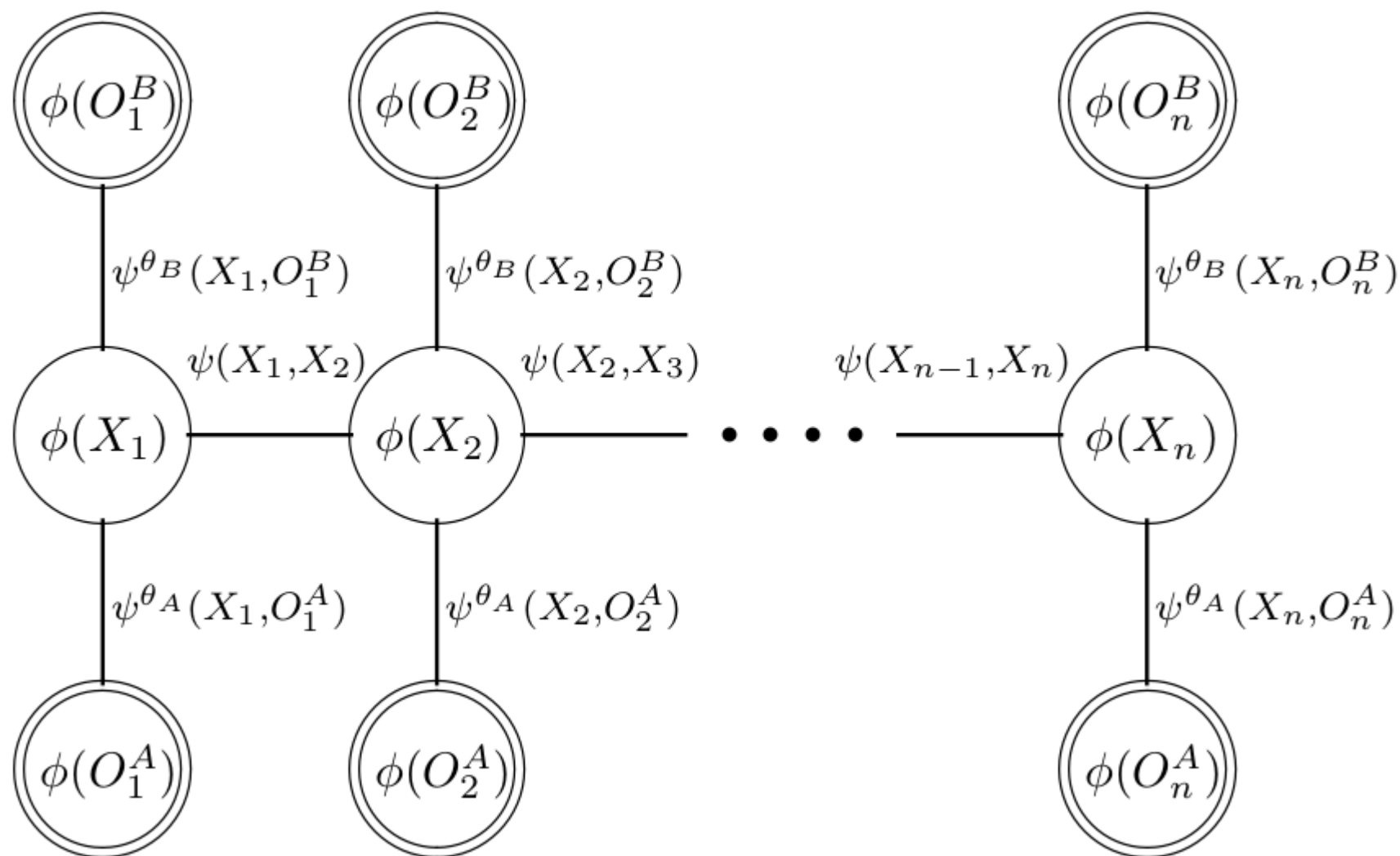
$$X_{AB} = \arg \min_X \theta_A E_A(X) + \theta_B E_B(X)$$

$$\text{s.t } \theta_A + \theta_B = 1$$

$$\theta_A, \theta_B \geq 0$$

$$\begin{aligned} X_{AB} &= \arg \min_X -\theta_A \left[ \sum_{i=1}^n \log \phi(X_i) + \sum_{i=1}^n \log \psi(\mathbf{O}_i^{\mathbf{A}}, X_i) + \sum_{i=1}^{n-1} \log \psi(X_i, X_{i+1}) \right] \\ &\quad - \theta_B \left[ \sum_{i=1}^n \log \phi(X_i) + \sum_{i=1}^n \log \psi(\mathbf{O}_i^{\mathbf{B}}, X_i) + \sum_{i=1}^{n-1} \log \psi(X_i, X_{i+1}) \right] \\ &= \arg \min_X - \left( \sum_{i=1}^n \log \phi(X_i) + \sum_{i=1}^{n-1} \log \psi(X_i, X_{i+1}) \right) \\ &\quad - \left( \sum_{i=1}^n (\log \psi(\mathbf{O}_i^{\mathbf{A}}, X_i)^{\theta_A} + \log \psi(\mathbf{O}_i^{\mathbf{B}}, X_i)^{\theta_B}) \right) \end{aligned}$$

# We get – Tree MRF(TMRFs)





The main steps in the max-product algorithm are :

*Calculation of the Optimal Probability value:*

$$V(i, k) = \begin{cases} \psi(\mathbf{O}_1^{\mathbf{A}}, X_1=k)^{\theta_A} \psi(\mathbf{O}_1^{\mathbf{B}}, X_1=k)^{\theta_B} \phi(X_1=k) & \text{if } i = 1 \\ \phi(X_i=k) \psi(\mathbf{O}_i^{\mathbf{A}}, X_i=k)^{\theta_A} \psi(\mathbf{O}_1^{\mathbf{B}}, X_1=k)^{\theta_B} \dots & \\ \max_j \{ \psi(X_i=k, X_{i-1}=j) V(i-1, j) \} & \text{if } 2 \leq i \leq n \end{cases}$$

*Retrieval of the optimal sequence corresponding to the optimal probability value:*

$$X_n^* = \arg \max_k V(n, k)$$

$$X_i^* = Ptr(X_{i+1}^*, i+1) \quad \text{if } 1 \leq i \leq n-1$$

---

**Algorithm 2** Pareto-Optimal MultiState Protein Design

---

```
1: procedure PARETO-FRONTIER( $\phi, \psi, O^A, O^B$ ) ▷ Find Pareto frontier
2:   Initialize queue  $Q \leftarrow \emptyset$ 
3:    $X_A \leftarrow \text{CMRF-DECODE}(\phi, \psi, O^A)$  ▷  $\arg \min_X E_A(X)$ 
4:    $X_B \leftarrow \text{CMRF-DECODE}(\phi, \psi, O^B)$  ▷  $\arg \min_X E_B(X)$ 
5:   Enqueue  $(X_A, X_B)$ 
6:    $C_H \leftarrow \{X_A, X_B\}$ 
7:   repeat
8:     Dequeue from Q a pair  $(X_1, X_2)$  and assert( $E_A(X_1) < E_A(X_2)$ )
9:      $m \leftarrow \frac{E_B(X_1) - E_B(X_2)}{E_A(X_1) - E_A(X_2)}$ 
10:     $\theta_A \leftarrow \frac{-m}{1-m}$ 
11:     $\theta_B \leftarrow \frac{1}{1-m}$ 
12:     $X_{AB} \leftarrow \text{TMRF-DECODE}(\phi, \psi, O^A, O^B, \theta_A, \theta_B)$ 
▷  $\arg \min_X \theta_A E_A(X) + \theta_B E_B(X)$ 
13:    if  $X_{AB} \neq X_A$  and  $X_{AB} \neq X_B$  then
14:       $C_H \leftarrow C_H \cup \{X_{AB}\}$ 
15:      Enqueue  $(X_A, X_{AB})$  and  $(X_{AB}, X_B)$ 
16:    end if
17:  until Q is empty
18:  return  $C_H$ 
19: end procedure
```

```

20: procedure CMRF-DECODE( $\phi, \psi, O^A$ ) ▷ Max decode CMRF
21:   Allocate  $V[N][K]$  ▷ N is length of seq, K is numAA
22:   Allocate  $X_A[N]$ 
23:   for  $k \leftarrow 1, K$  do
24:      $V[1][k] \leftarrow \psi(\mathbf{O}_1^A, X_1=k)\phi(X_1=k)$ 
25:   end for
26:   for  $i \leftarrow 1, N$  do
27:     for  $k \leftarrow 1, K$  do
28:        $V[i][k] \leftarrow \phi(X_i=k)\psi(\mathbf{O}_i^A, X_i=k) \max_j \{\psi(X_i=k, X_{i-1}=j)V(i -$ 
29:          $1, j)\}$ 
30:        $Ptr[i][k] \leftarrow \arg \max_j \{\psi(X_i=k, X_{i-1}=j)V(i - 1, j)\}$ 
31:     end for
32:   end for ▷ Now retrieve the pointers
33:    $X_A[N] \leftarrow \arg \max_k V[N][k]$ 
34:   for  $i \leftarrow N - 1, 1$  do
35:      $X_A[i] \leftarrow Ptr[i + 1][X_A[i + 1]]$ 
36:   end for
37:   return  $X_A$ 
end procedure

```

# Summary

## **Contributions**

A fast and exact protein design method

Flexibility in choosing features

Enumerate the Pareto Frontier

## **Limitations**

Sequential dependency between variables

Thank You