
Modeling and predicting the structure of transmembrane proteins

Jérôme Waldispühl ¹²³, Jean-Marc Steyaert ²

`Jerome.Waldispuhl@polytechnique.edu`

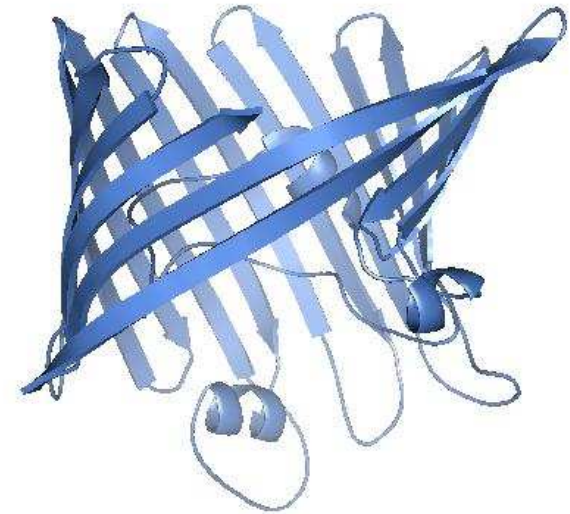
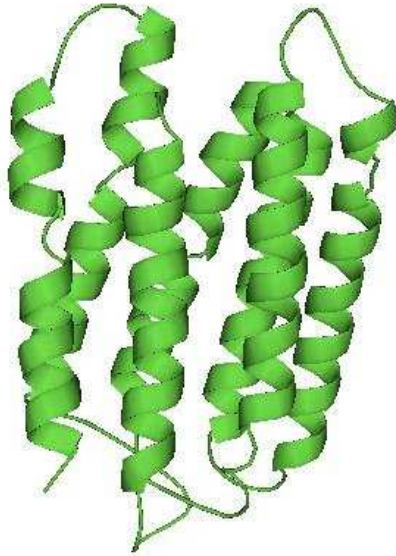
¹ Department of Biology, Boston College, USA

² LIX, École polytechnique, France

³ LIAFA, Université de Paris 7, France



Modeling and predicting the structure of transmembrane proteins



- Summary of language theory (multi-tape S-attribute grammars),
- some notions of biology,
- an approximate physical model,
- grammatical modeling,
- performance evaluation,
- conclusion.

Definition *Context-free grammars*

$$G = \{V_T, V_N, P, S\}$$

- V_T is the set of terminals,
- V_N is the set of non-terminals,
- P is the set of productions rules ($A \rightarrow \alpha$),
- S is the axiom,

Definition *S-attribute grammars*

$$G = \{V_T, V_N, P, S, \mathcal{A}, \lambda_{\mathcal{A}}, F_P\}$$

- V_T is the set of terminals,
- V_N is the set of non-terminals,
- P is the set of productions rules,
- S is the axiom,
- \mathcal{A} is the set of attributes,
- $\lambda_{\mathcal{A}}$ is the set of evaluation functions for the terminals,
- F_P is the set of the functions used to compute the non-terminals attributes.

Example (1) : arithmetical expressions

$$V_T = \{0, \dots, 9, +, \times\},$$

$$V_N = S,$$

$$P = \left\{ \begin{array}{l} S \rightarrow S + S \\ S \rightarrow S \times S \\ S \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{array} \right.$$

Context-free grammar

Example (1) : arithmetical expressions

$$\left. \begin{array}{l} V_T = \{0, \dots, 9, +, \times\}, \\ V_N = S, \\ P = \begin{cases} S \rightarrow S + S \\ S \rightarrow S \times S \\ S \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{cases} \end{array} \right) \text{Context-free grammar}$$

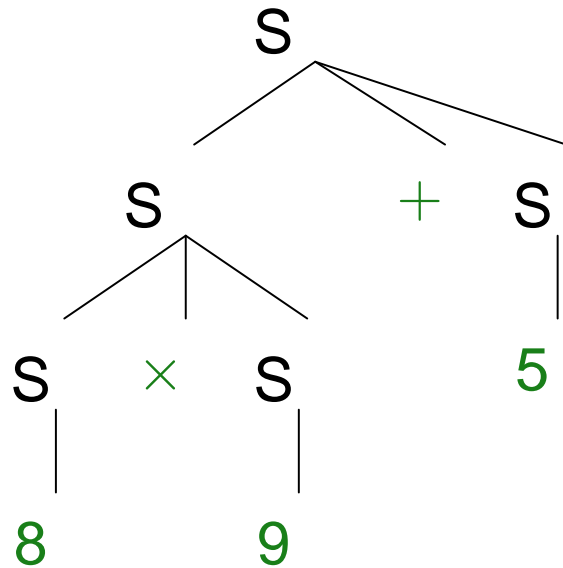
$$\left. \begin{array}{l} \mathcal{A} = \mathbb{N} \\ \lambda_{\mathcal{A}} = \begin{cases} S_{\mathcal{A}}(0) = 0 \\ \vdots \\ S_{\mathcal{A}}(9) = 9 \\ S_{\mathcal{A}}(+, \times) = 0 \end{cases} \\ F_P = \begin{cases} f_{S \rightarrow S+S}(xyz) = x + z \\ f_{S \rightarrow S \times S}(xyz) = x \times z \\ f_{S \rightarrow a \in V_T}(x) = x \end{cases} \end{array} \right) \text{Attribute system}$$

Exemple (1) : Arithmetical expressions

$$8 \times 9 + 5$$

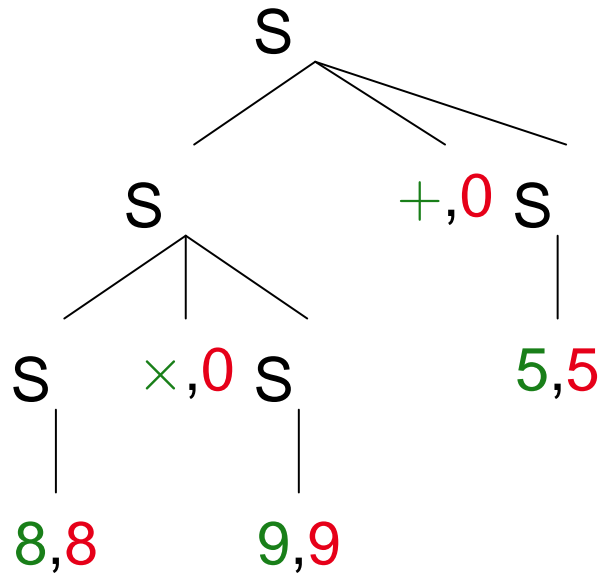
Exemple (1) : Arithmetical expressions

$$8 \times 9 + 5$$



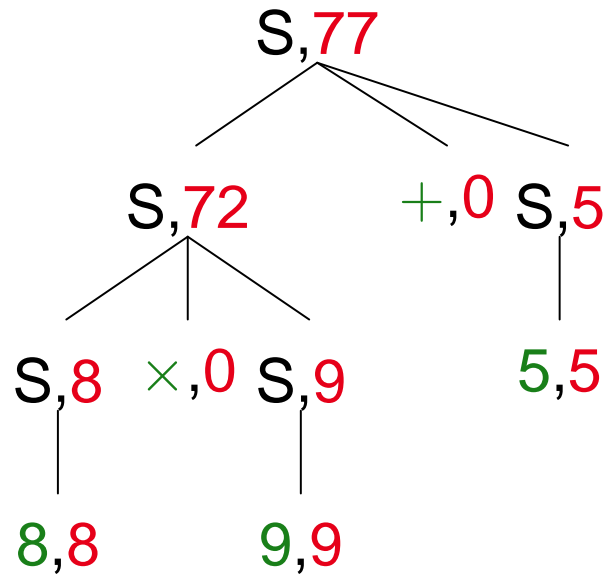
Exemple (1) : Arithmetical expressions

$$8 \times 9 + 5$$



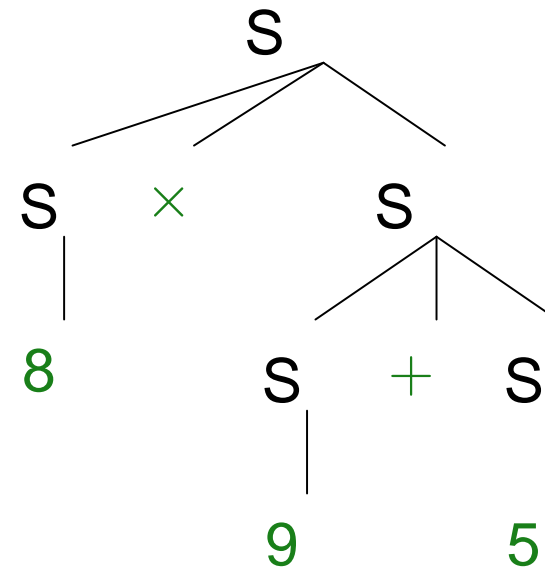
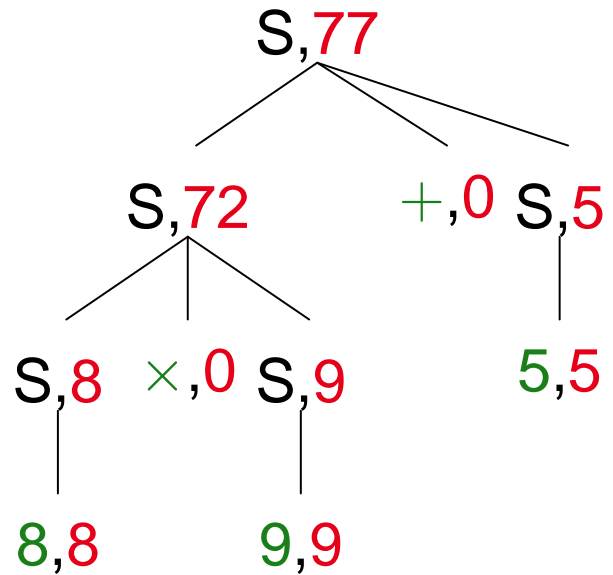
Exemple (1) : Arithmetical expressions

$$8 \times 9 + 5$$



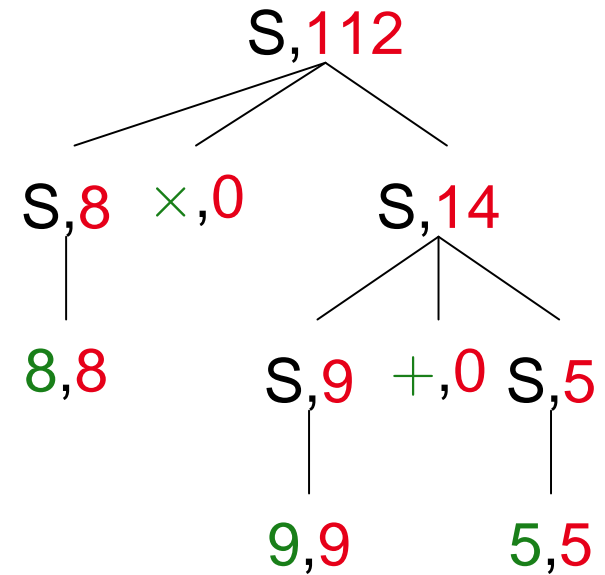
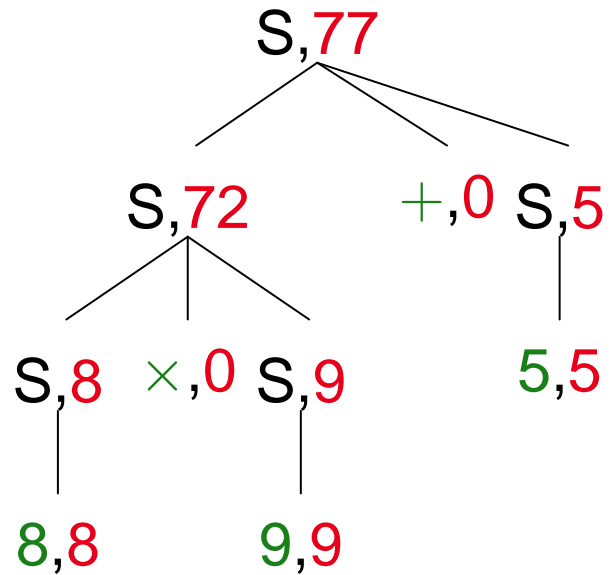
Exemple (1) : Arithmetical expressions

$$8 \times 9 + 5$$

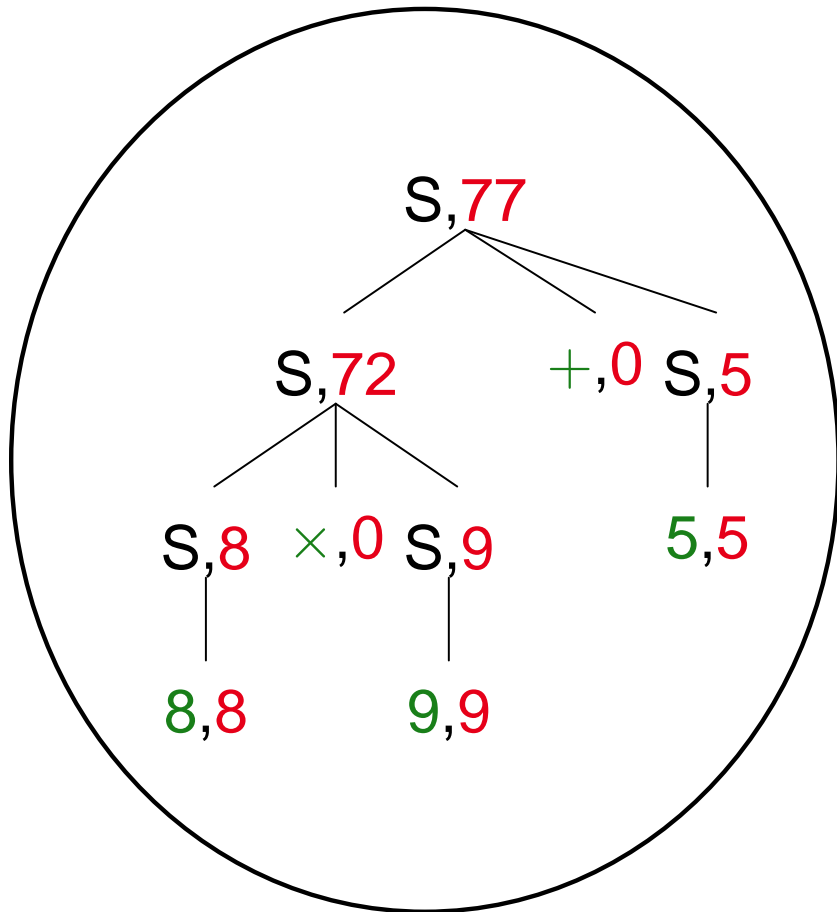


Exemple (1) : Arithmetical expressions

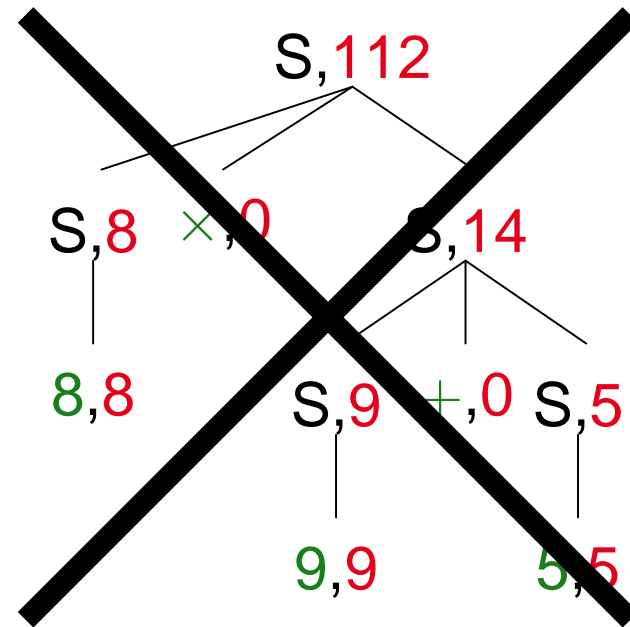
$$8 \times 9 + 5$$



Exemple (1) : Arithmetical expressions



$8 \times 9 + 5$



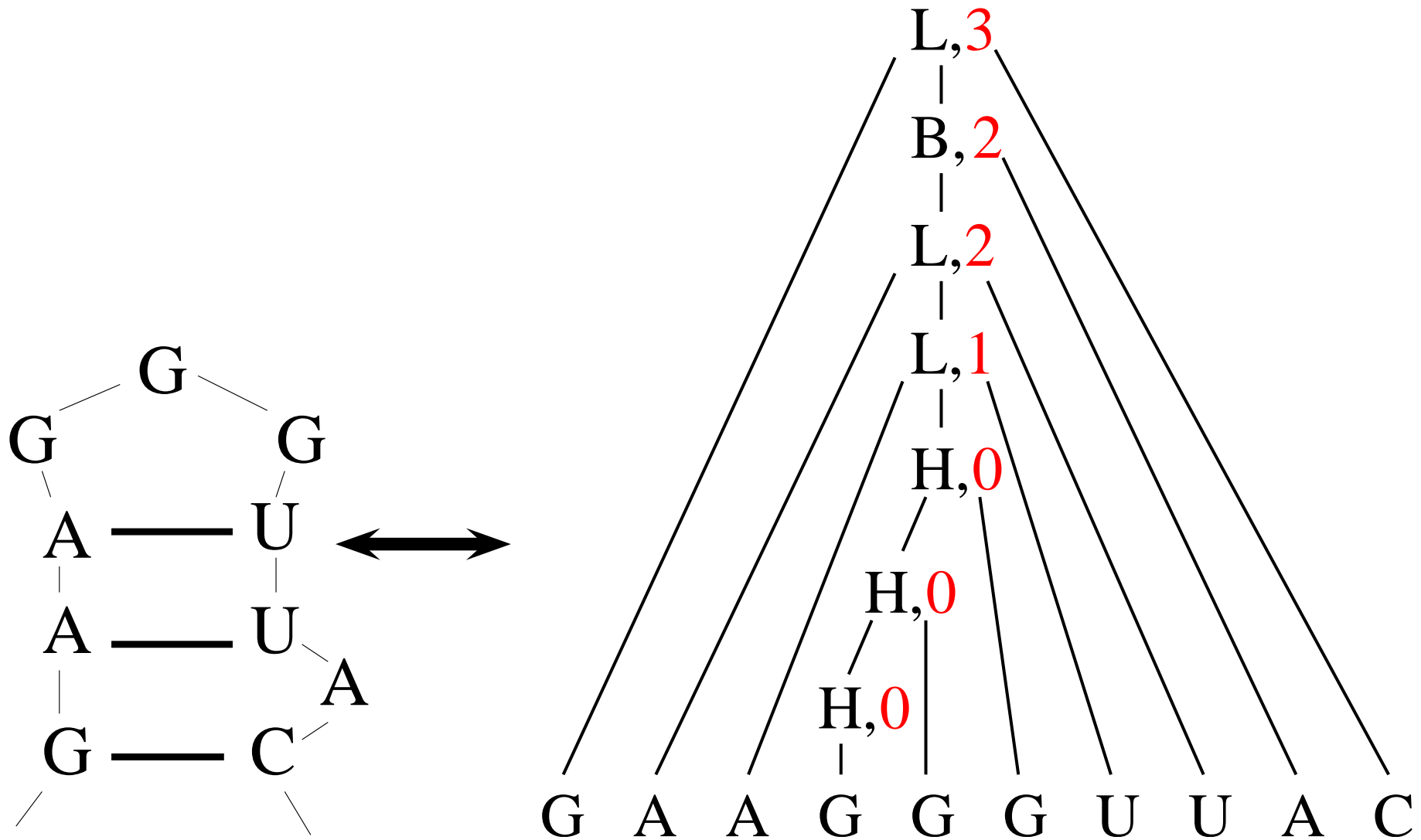
...Only the derivation tree with the optimal attribute is conserved.

Definition *Optimization constraint*

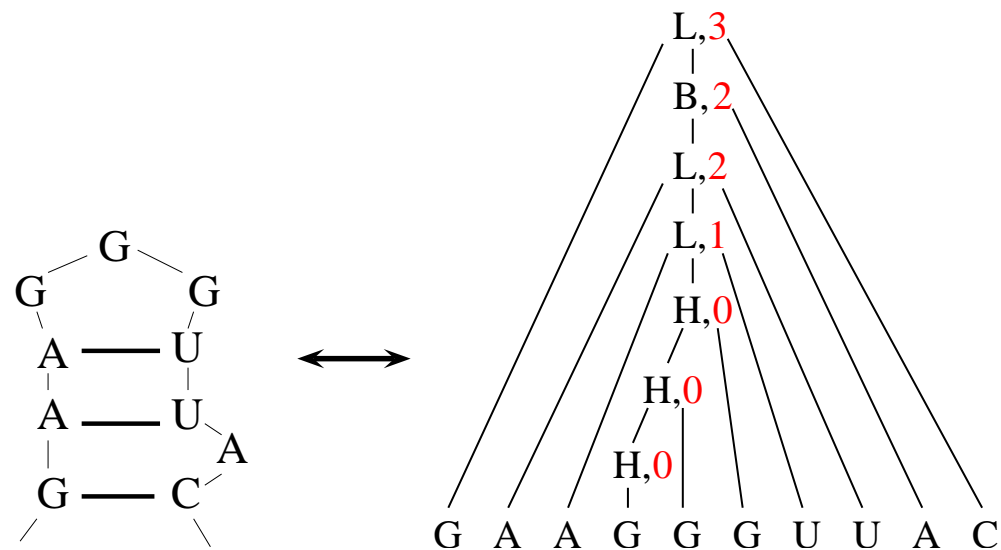
$$\mathcal{C}(x, \lambda_x, y, \lambda_y)$$

- x and $y \in V_T \cup V_N$,
- λ_x and $\lambda_y \in \mathcal{A}$,
- return the pair (z, λ_z) such that λ_z is optimal.

Example (2) : RNA secondary structure



Example (2) : RNA secondary structure



- secondary structure = **derivation tree**
- folding energy = **attributes**
- secondary structure with the minimum free energy (Zuker) = **derivation tree with the optimal attribute**

How to find the optimal derivation tree ?

principle : dynamic programming

algorithm : Cocke-Kazamy-Younger, Earley, GCP...

implementation : *mtsag2c* (F. Lefebvre, 1997)

Multi-tape S-attribute grammar

Definition *Multi-tape alphabet*

$$\Sigma = \bigotimes_{i=1 \dots m} (\Sigma^{(i)} \cup \{\varepsilon\})$$

Definition *Multi-tape Context-free grammar*

$$G = \{V_T, V_N, P, S\} \text{ where } V_T \text{ is an } m\text{-tape alphabet.}$$

Definition *Multi-tape S-attribute grammar*

$$G = \{V_T, V_N, P, S, \mathcal{A}, \lambda_{\mathcal{A}}, F_P\} \text{ where } V_T \text{ is an } m\text{-tape alphabet.}$$

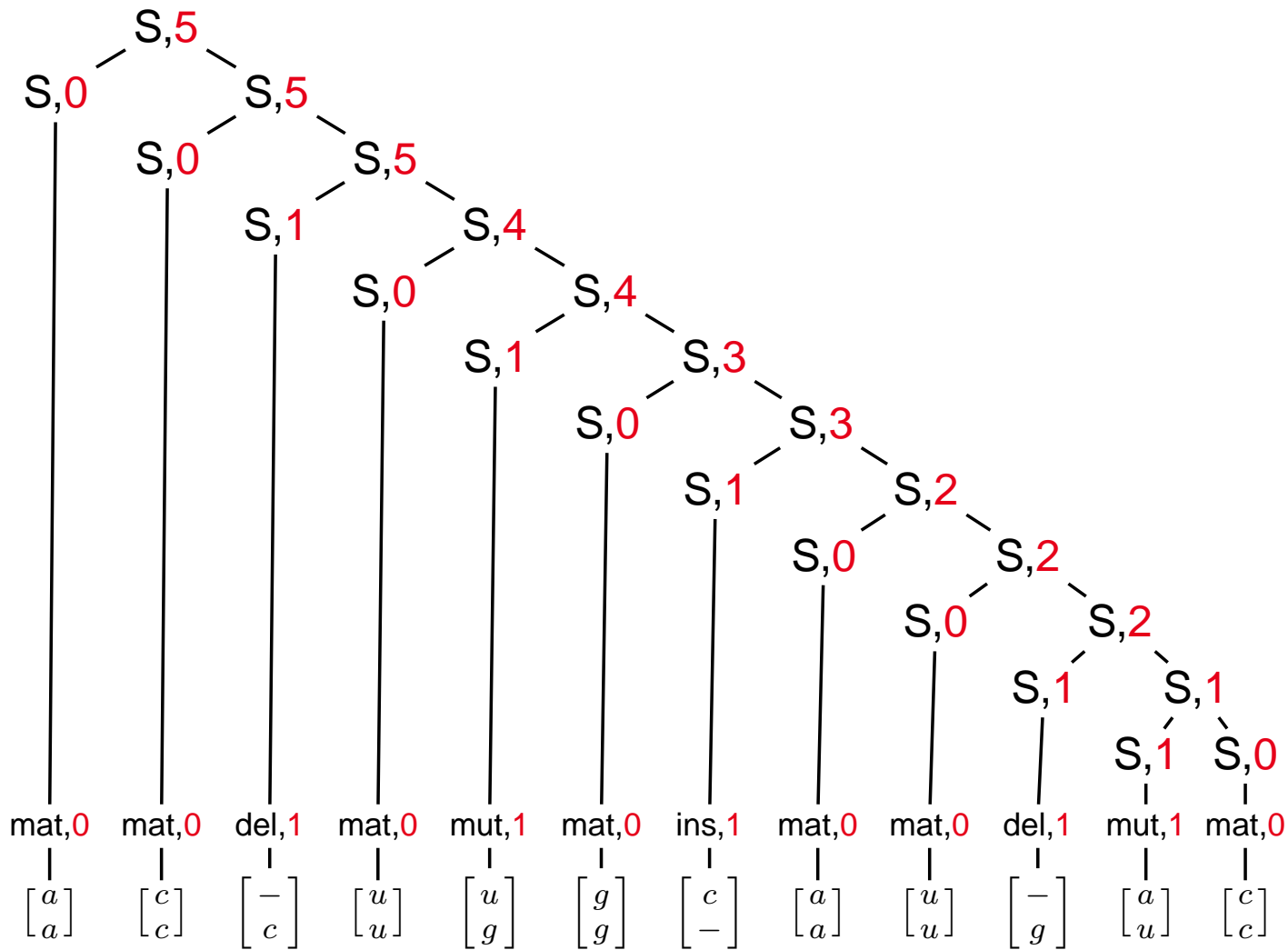
Example : RNA sequence alignment

$$\left\{ \begin{array}{l} S \rightarrow SS \mid mat \mid del \mid ins \mid mut \\ mat \rightarrow \begin{bmatrix} a \\ a \end{bmatrix} \mid \begin{bmatrix} u \\ u \end{bmatrix} \mid \begin{bmatrix} g \\ g \end{bmatrix} \mid \begin{bmatrix} c \\ c \end{bmatrix} \\ del \rightarrow \begin{bmatrix} - \\ a \end{bmatrix} \mid \begin{bmatrix} - \\ u \end{bmatrix} \mid \begin{bmatrix} - \\ g \end{bmatrix} \mid \begin{bmatrix} - \\ c \end{bmatrix} \\ ins \rightarrow \begin{bmatrix} a \\ - \end{bmatrix} \mid \begin{bmatrix} u \\ - \end{bmatrix} \mid \begin{bmatrix} g \\ - \end{bmatrix} \mid \begin{bmatrix} c \\ - \end{bmatrix} \\ mut \rightarrow \begin{bmatrix} a \\ u \end{bmatrix} \mid \begin{bmatrix} a \\ g \end{bmatrix} \mid \begin{bmatrix} a \\ c \end{bmatrix} \mid \begin{bmatrix} u \\ a \end{bmatrix} \\ \quad \mid \begin{bmatrix} u \\ g \end{bmatrix} \mid \begin{bmatrix} u \\ c \end{bmatrix} \mid \begin{bmatrix} g \\ a \end{bmatrix} \mid \begin{bmatrix} g \\ u \end{bmatrix} \\ \quad \mid \begin{bmatrix} g \\ c \end{bmatrix} \mid \begin{bmatrix} c \\ a \end{bmatrix} \mid \begin{bmatrix} c \\ u \end{bmatrix} \mid \begin{bmatrix} c \\ g \end{bmatrix} \end{array} \right.$$

Example : RNA sequence alignment

$$\left\{ \begin{array}{l}
 S \rightarrow SS \mid mat \mid del \mid ins \mid mut \\
 mat \rightarrow \begin{bmatrix} a \\ a \end{bmatrix} \mid \begin{bmatrix} u \\ u \end{bmatrix} \mid \begin{bmatrix} g \\ g \end{bmatrix} \mid \begin{bmatrix} c \\ c \end{bmatrix} \\
 del \rightarrow \begin{bmatrix} - \\ a \end{bmatrix} \mid \begin{bmatrix} - \\ u \end{bmatrix} \mid \begin{bmatrix} - \\ g \end{bmatrix} \mid \begin{bmatrix} - \\ c \end{bmatrix} \\
 ins \rightarrow \begin{bmatrix} a \\ - \end{bmatrix} \mid \begin{bmatrix} u \\ - \end{bmatrix} \mid \begin{bmatrix} g \\ - \end{bmatrix} \mid \begin{bmatrix} c \\ - \end{bmatrix} \\
 mut \rightarrow \begin{bmatrix} a \\ u \end{bmatrix} \mid \begin{bmatrix} a \\ g \end{bmatrix} \mid \begin{bmatrix} a \\ c \end{bmatrix} \mid \begin{bmatrix} u \\ a \end{bmatrix} \\
 \quad \mid \begin{bmatrix} u \\ g \end{bmatrix} \mid \begin{bmatrix} u \\ c \end{bmatrix} \mid \begin{bmatrix} g \\ a \end{bmatrix} \mid \begin{bmatrix} g \\ u \end{bmatrix} \\
 \quad \mid \begin{bmatrix} g \\ c \end{bmatrix} \mid \begin{bmatrix} c \\ a \end{bmatrix} \mid \begin{bmatrix} c \\ u \end{bmatrix} \mid \begin{bmatrix} c \\ g \end{bmatrix}
 \end{array} \right. \quad F_P = \left\{ \begin{array}{l}
 \mathcal{A} = \mathbb{Z} \\
 \lambda(\bullet) = 0 \\
 f_{S \rightarrow SS}(xy) = x + y \\
 f_{S \rightarrow del \mid ins \mid mut}(x) = x \\
 f_{mat \rightarrow \bullet}(x) = 0 \\
 f_{del \rightarrow \bullet}(x) = 1 \\
 f_{ins \rightarrow \bullet}(x) = 1 \\
 f_{mut \rightarrow \bullet}(x) = 1
 \end{array} \right.$$

Example : RNA sequence alignment

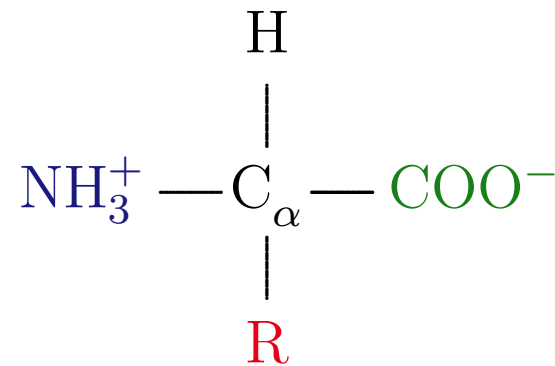


Notions of Biology

- Definition of a protein
- Structure of proteins
- Transmembrane channels

Definition of a protein

Amino acid chemical formula :



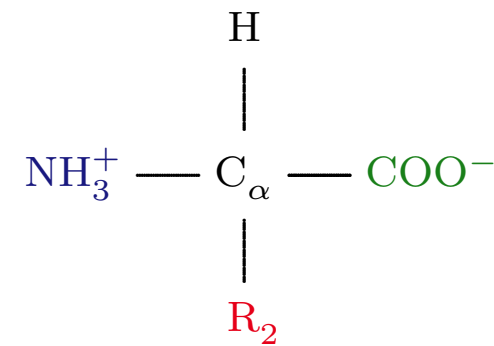
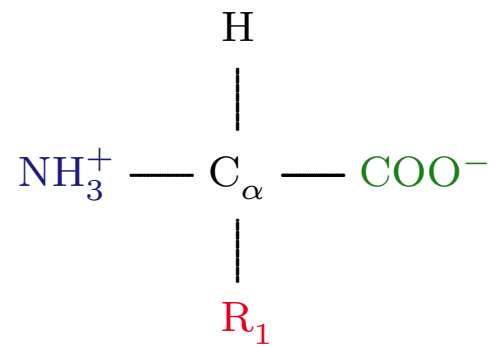
Definition of a protein

The 20 amino acids :

$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{H} \end{array}$ <p><i>Glycine (Gly/G)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH}_3 \end{array}$ <p><i>Alanine (Ala/A)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH} \\ / \quad \backslash \\ \text{CH}_3 \quad \text{CH}_3 \end{array}$ <p><i>Valine (Val/V)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH}_2 \\ \\ \text{CH} \\ / \quad \backslash \\ \text{CH}_3 \quad \text{CH}_3 \end{array}$ <p><i>Leucine (Leu/L)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{H} - \text{C} - \text{CH}_3 \\ \\ \text{CH}_2 \\ \\ \text{CH}_3 \end{array}$ <p><i>Isoleucine (Ile/L)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH}_2 \\ \\ \text{OH} \end{array}$ <p><i>Serine (Ser/S)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{H} - \text{C} - \text{OH} \\ \\ \text{CH}_3 \end{array}$ <p><i>Threonine (Thr/T)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH}_2 \\ \\ \text{C}_6\text{H}_5 \end{array}$ <p><i>Phenylalanine (Phe/F)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH}_2 \\ \\ \text{C}_6\text{H}_4 \\ \\ \text{OH} \end{array}$ <p><i>Tyrosine (Tyr/Y)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH}_2 \\ \\ \text{C} = \text{CH} \\ \quad \backslash \\ \text{C}_6\text{H}_4 \quad \text{NH} \end{array}$ <p><i>Tryptophane (Trp/W)</i></p>
$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{NH}_3^+ \end{array}$ <p><i>Lysine (Lys/K)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{C} = \text{NH}_2 \\ \\ \text{NH}_2 \end{array}$ <p><i>Arginine (Arg/R)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH}_2 \\ \\ \text{C} \\ / \quad \backslash \\ \text{CH} \quad \text{NH} \\ \quad \backslash \\ \text{N} \quad \text{CH} \end{array}$ <p><i>Histidine (His/H)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH}_2 \\ \\ \text{COO}^- \end{array}$ <p><i>Aspartate (Asp/D)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{COO}^- \end{array}$ <p><i>Glutamate (Glu/E)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{NH}_2 \end{array}$ <p><i>Asparagine (Asn/N)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{NH}_2 \end{array}$ <p><i>Glutamamine (Gln/Q)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH}_2 \\ \\ \text{SH} \end{array}$ <p><i>Cystéine (Cys/C)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{S} \\ \\ \text{CH}_3 \end{array}$ <p><i>Méthionine (Met/M)</i></p>	$\begin{array}{c} \text{H} \\ \\ \text{NH}_3^+ - \text{C}_\alpha - \text{COO}^- \\ \\ \text{CH}_2 \\ / \quad \backslash \\ \text{CH}_2 \quad \text{CH}_2 \end{array}$ <p><i>Proline (Pro/P)</i></p>

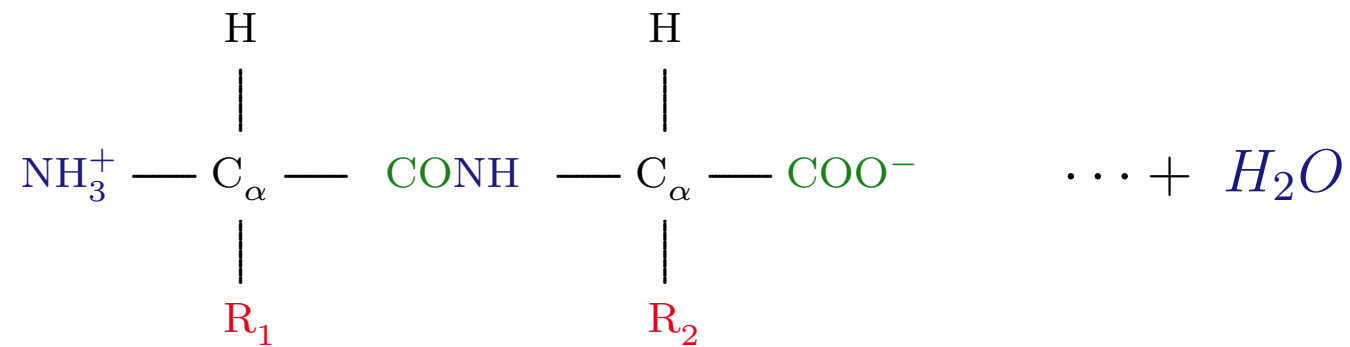
Definition of a protein

The peptid bond :

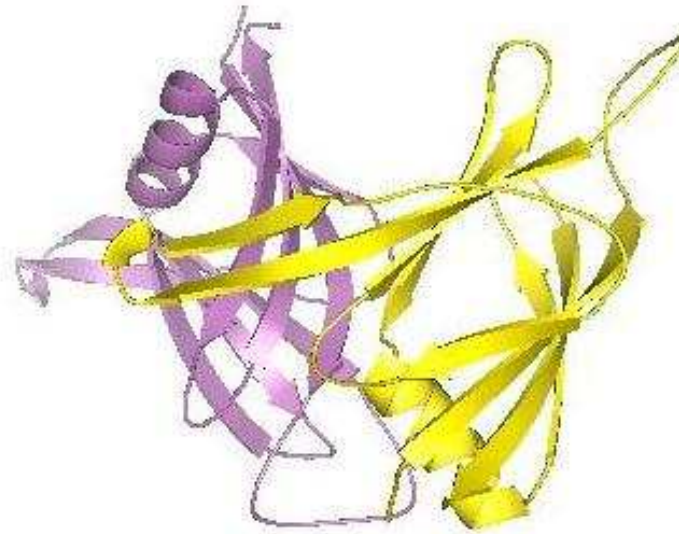


Definition of a protein

The peptid bond :

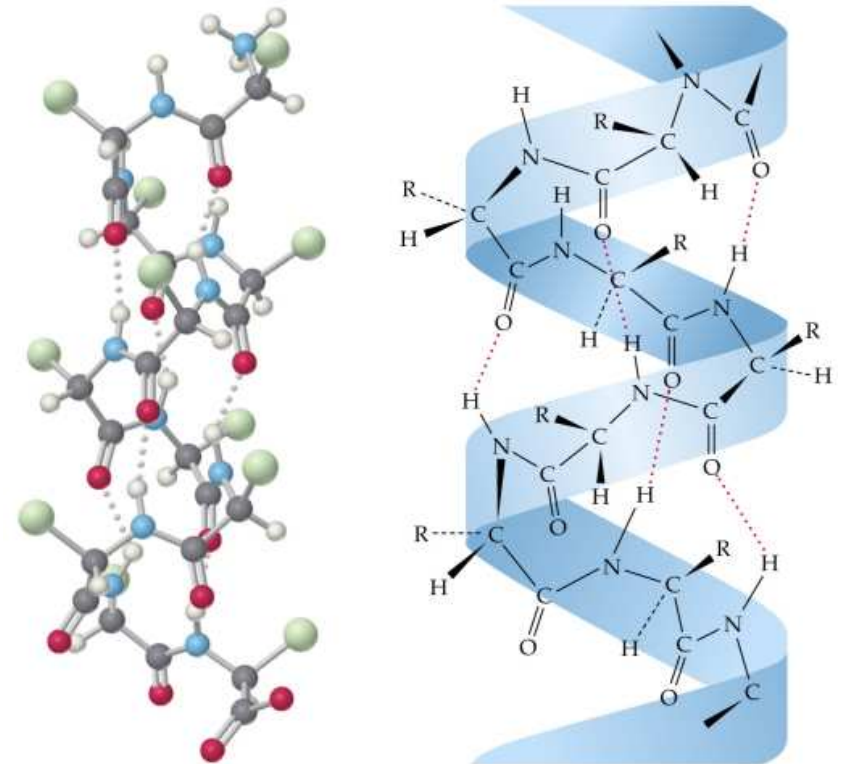


Kynase C



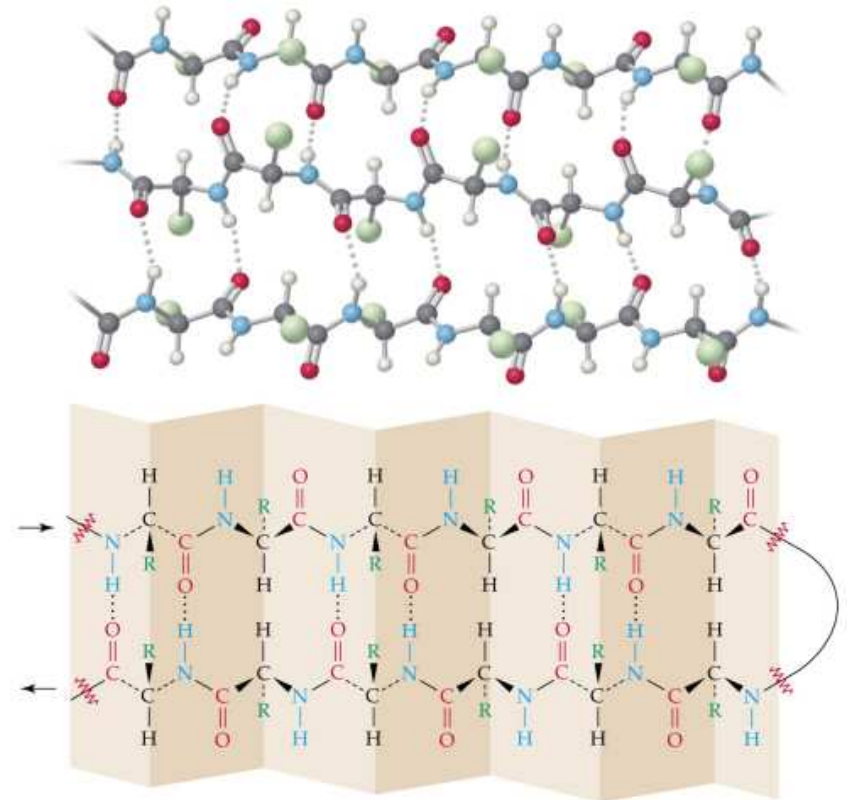
α -helix

- 3.6 amino acids per turn,
- hydrogen bond between residus n and $n + 4$.



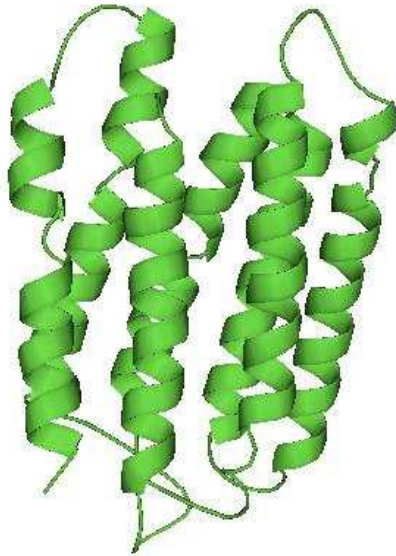
β -sheet

- composed of β -strands
- 2 amino acids per turn,
- hydrogen bond between residues of paired β -strands.

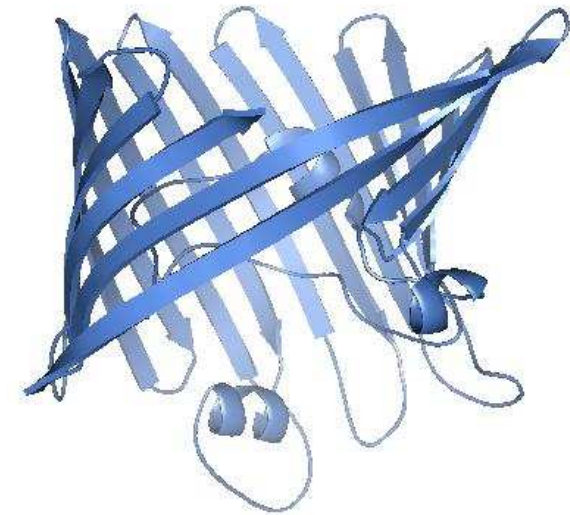


Transmembrane channels

Bacteriorhodopsin

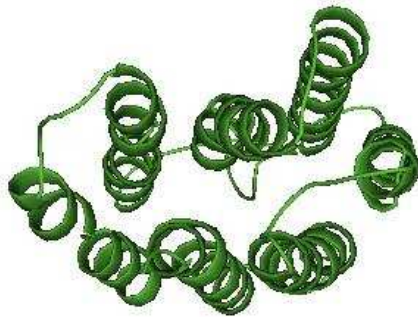


Porin



Transmembrane channels

Bacteriorhodopsin



Porin



Transmembrane channel

Why ?

Transmembrane channel

Why ?

- Simple topologies (only parallel or anti-parallel pairings),
- strong constraints from the environment,
- Some parameters are (much) more important than the others (hydrophobicity)

Transmembrane channel

Why ?

- Simple topologies (only parallel or anti-parallel pairings),
- strong constraints from the environment,
- Some parameters are (much) more important than the others (hydrophobicity)

Interest ?

Transmembrane channel

Why ?

- Simple topologies (only parallel or anti-parallel pairings),
- strong constraints from the environment,
- Some parameters are (much) more important than the others (hydrophobicity)

Interest ?

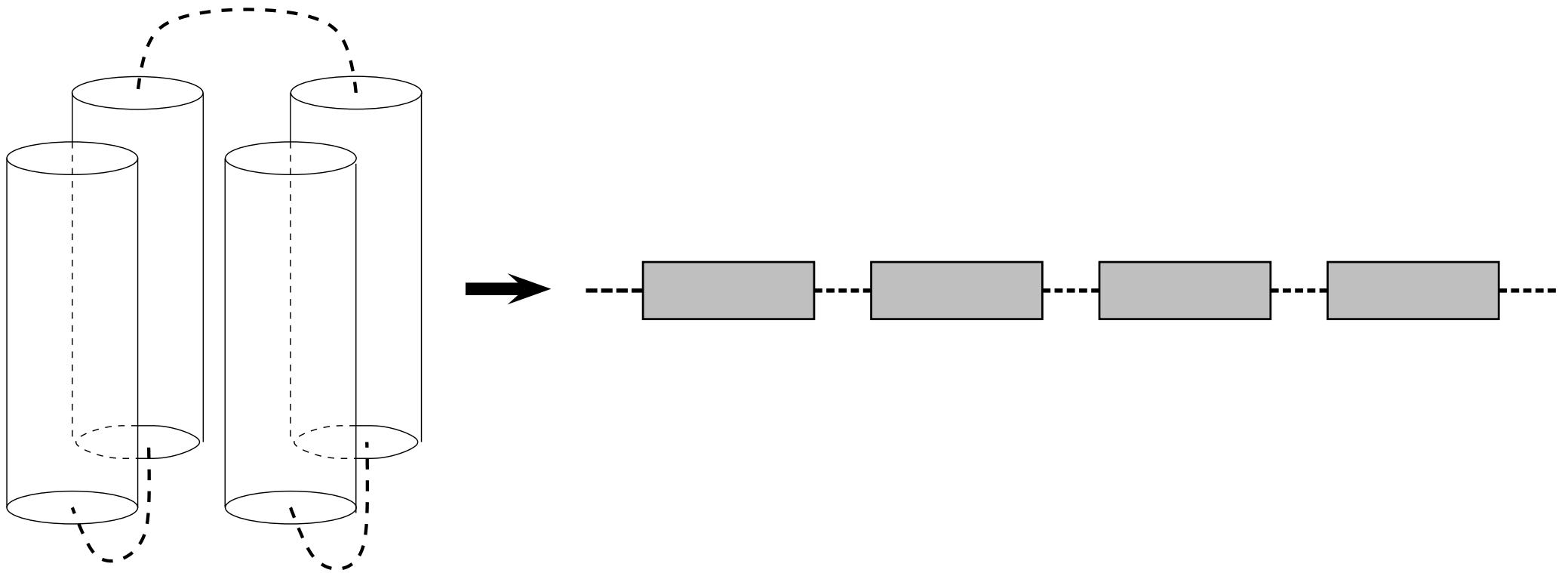
- nearly 40% of the proteome,
- functional importance (allows communication between inner and outer milieu of cell),
- difficult to be observe experimentally.

Approximate physical model for transmembrane channels

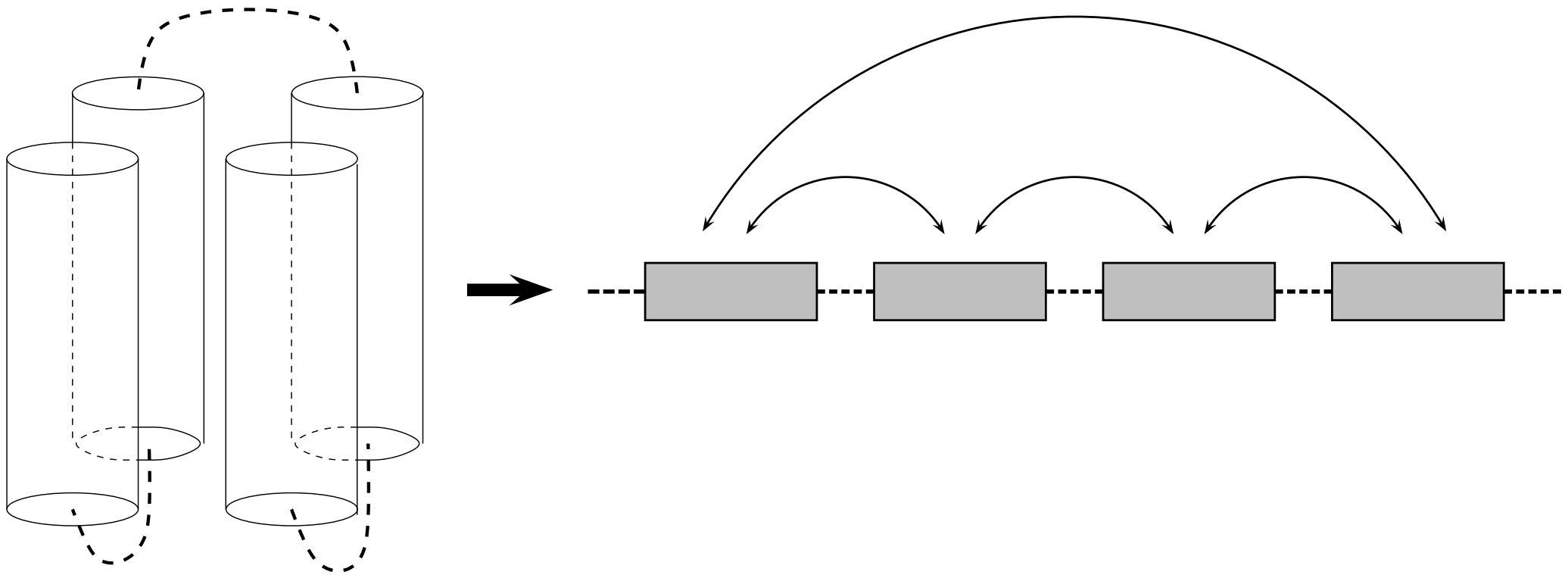
Approximate physical model for α -transmembrane channels

- Modeling the overall structure of α -channel,
- modeling anti-parallel pairing of α -helices,
- modeling the local structure of α -helices,
- pseudo folding energy of α -channels.

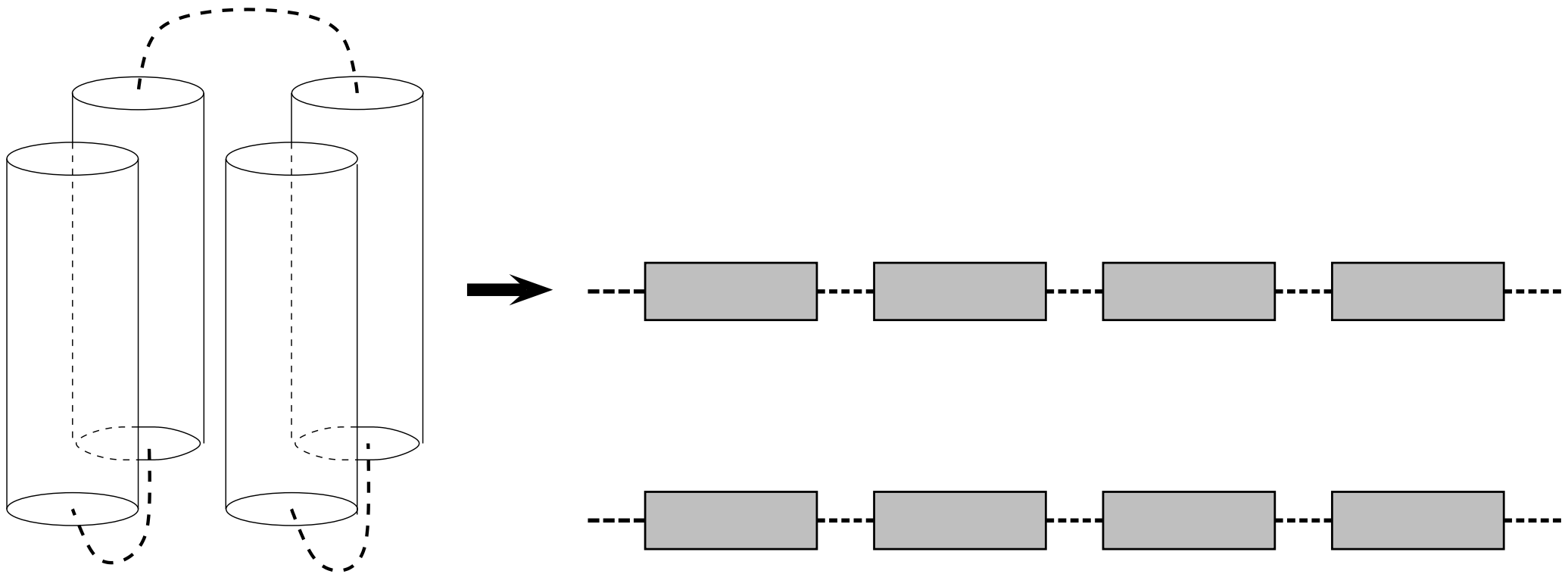
Modeling the overall structure of α -channel



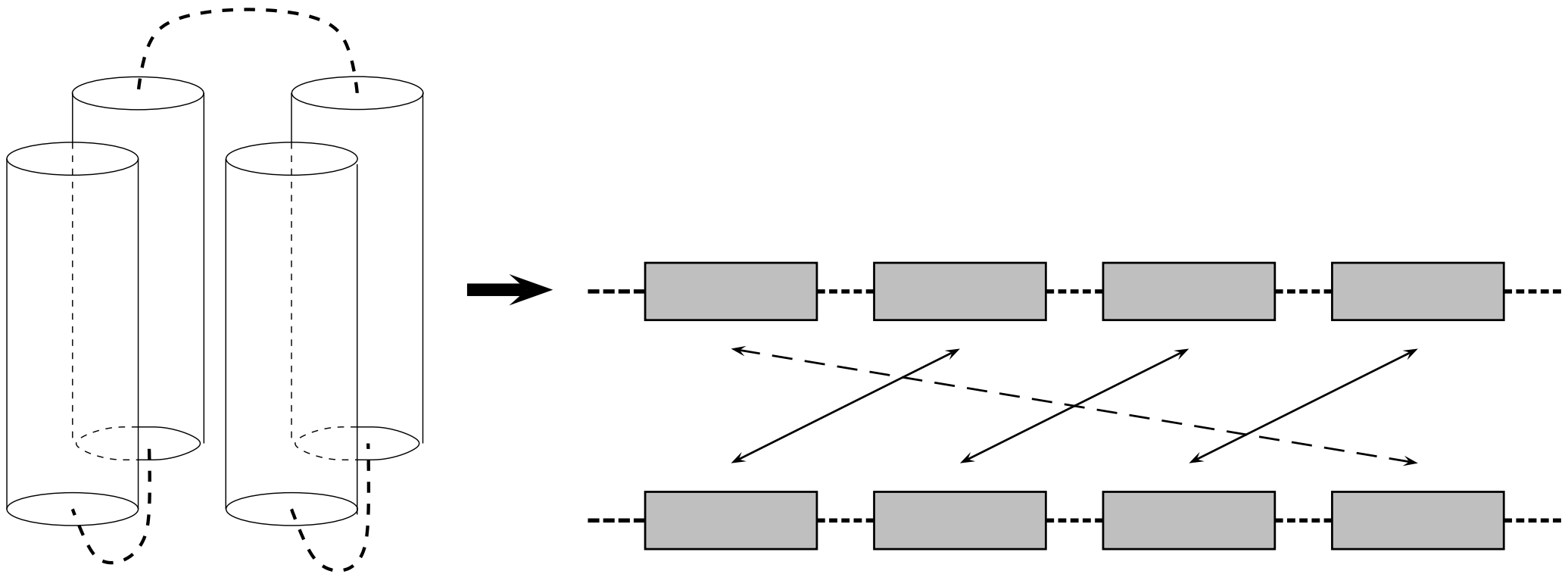
Modeling the overall structure of α -channel



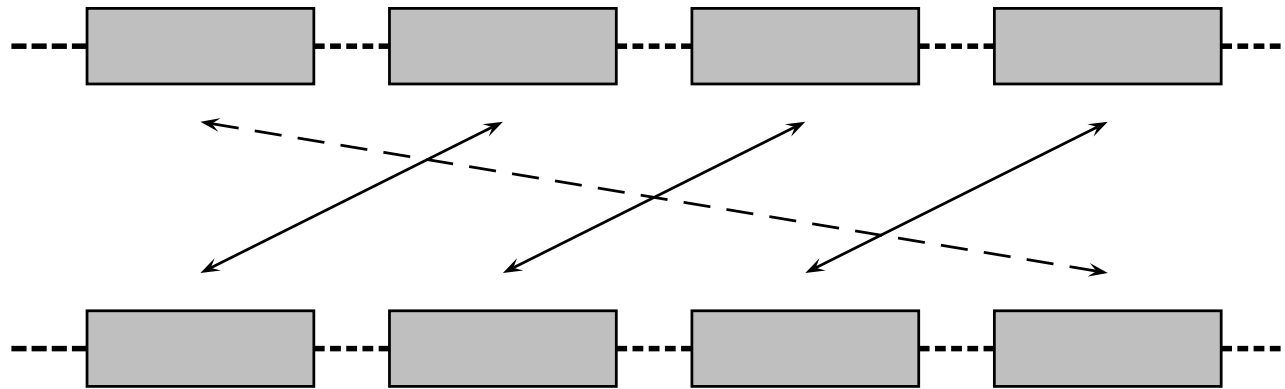
Modeling the overall structure of α -channel



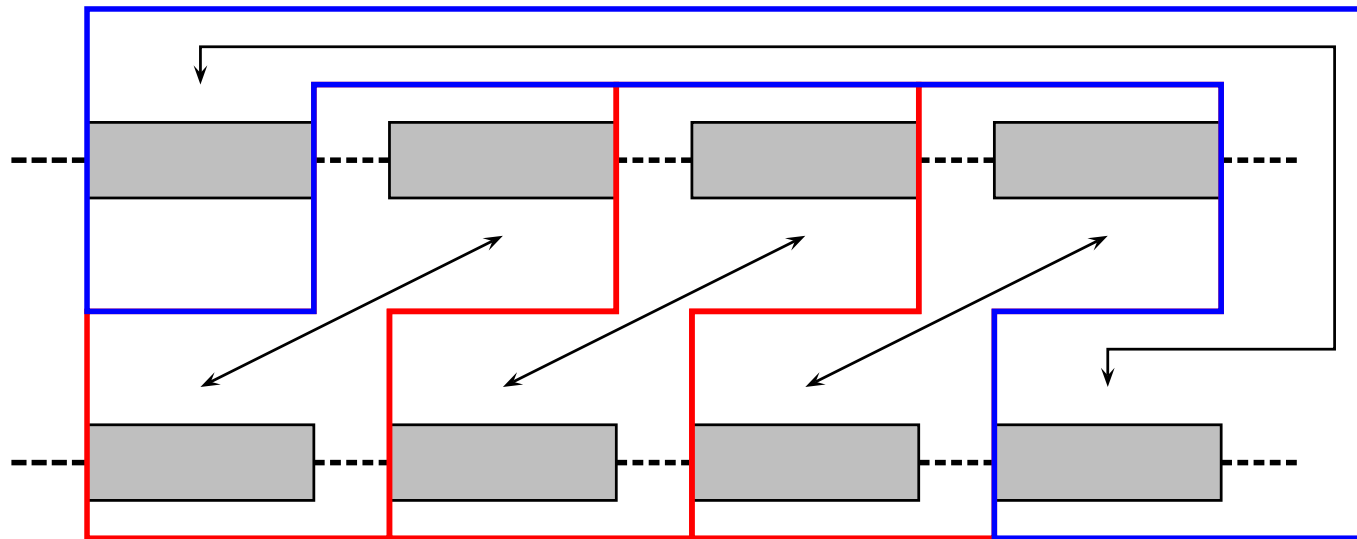
Modeling the overall structure of α -channel



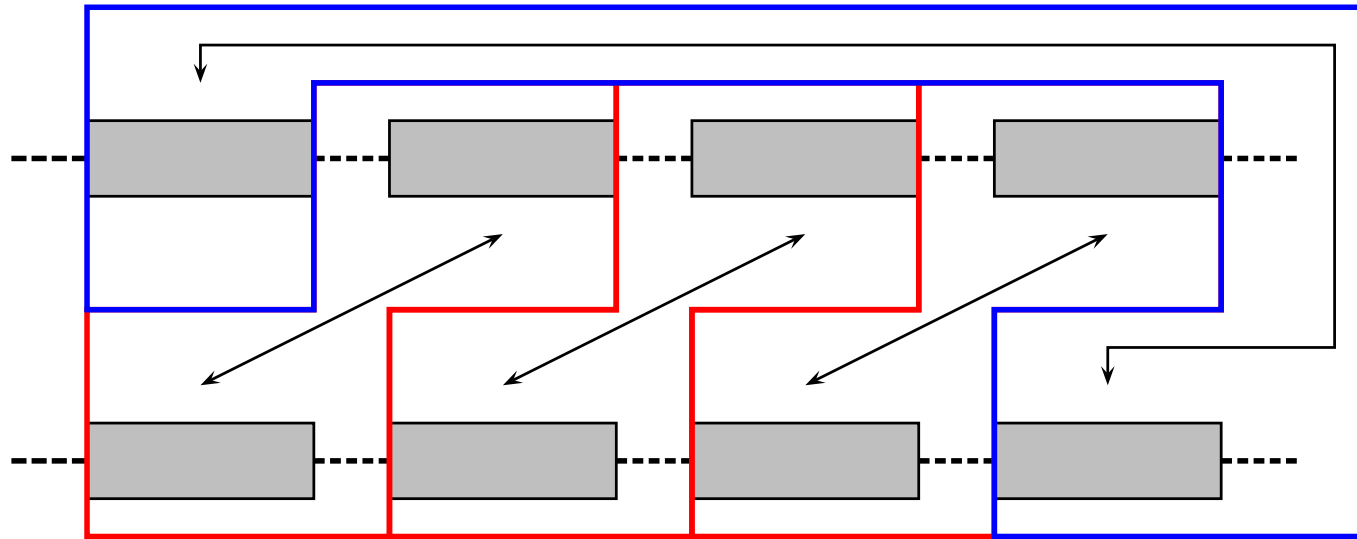
Modeling the overall structure of α -channel



Modeling the overall structure of α -channel

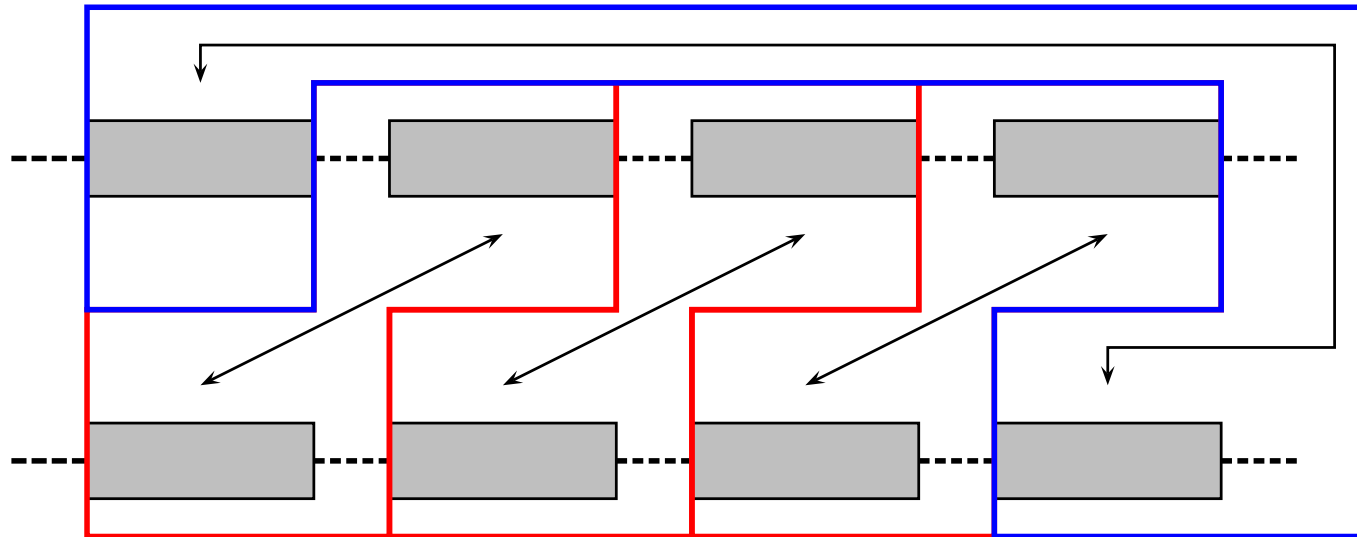


Modeling the overall structure of α -channel



Description of α -channels with *only* simple anti-parallel pairings.

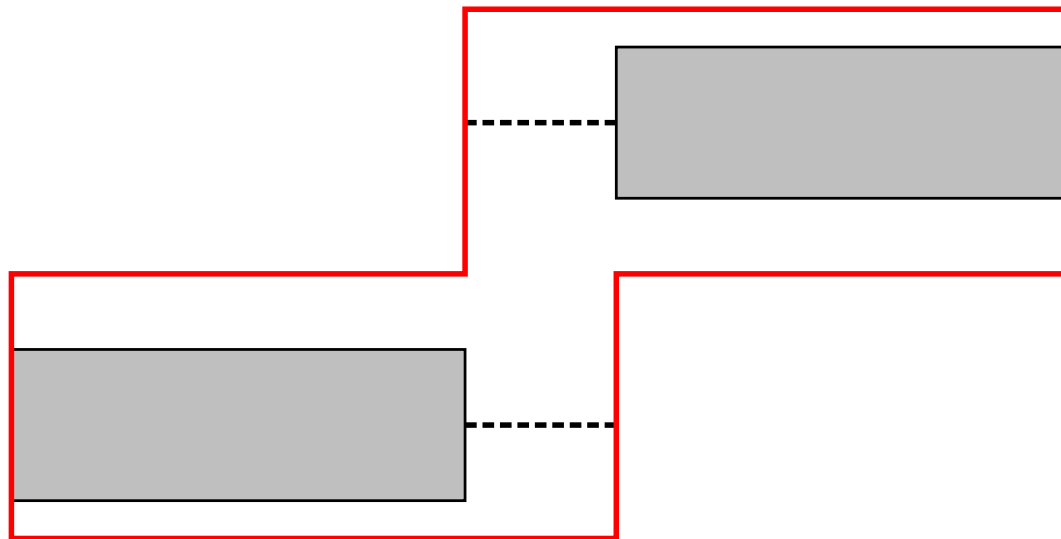
Modeling the overall structure of α -channel



An α -channel is a concatenation of simple anti-parallel pairings.

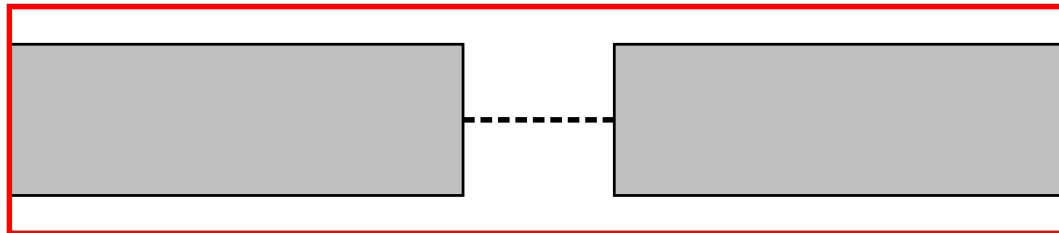
Modeling anti-parallel pairing of α -helices

Let's go back to a linear description :

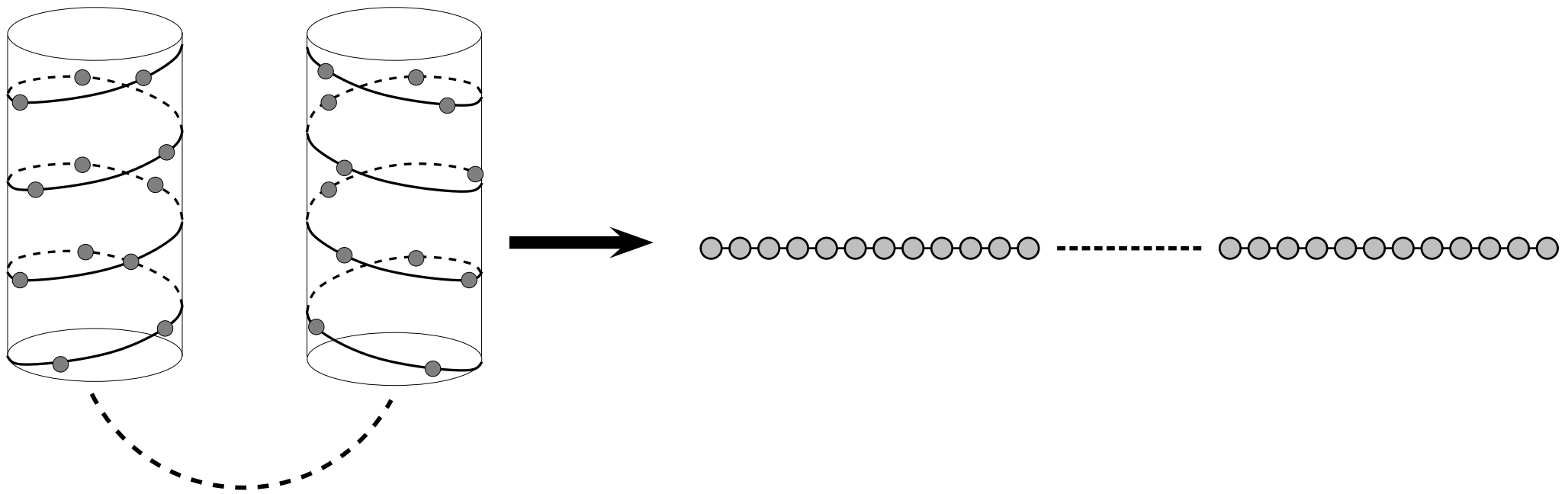


Modeling anti-parallel pairing of α -helices

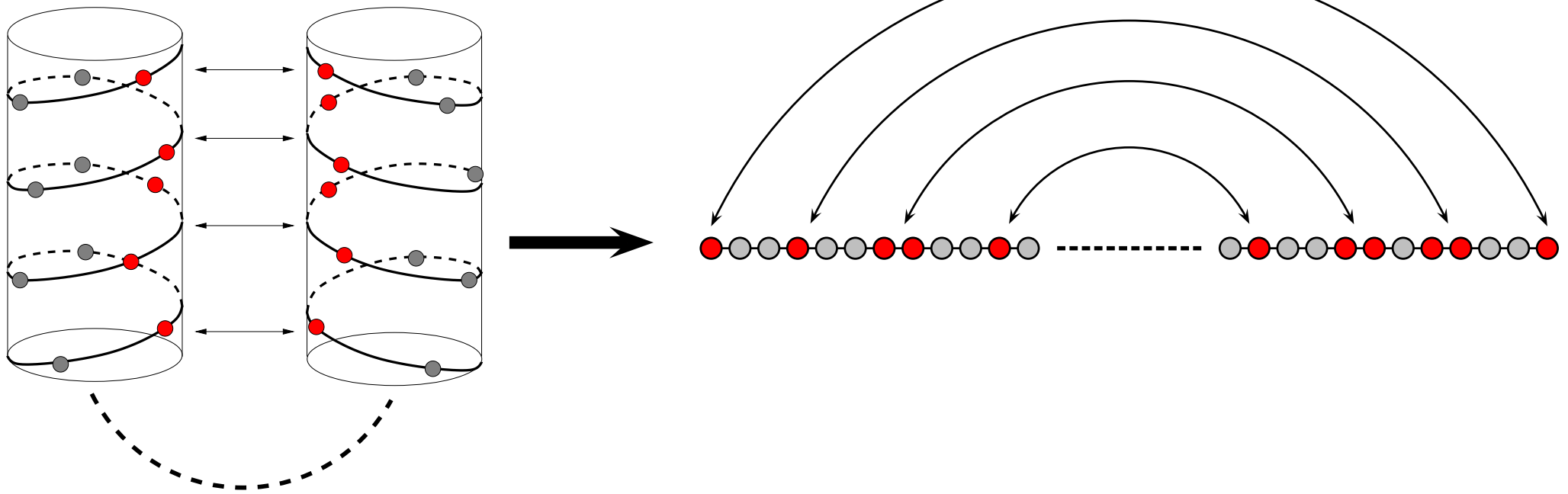
Let's go back to a linear description :



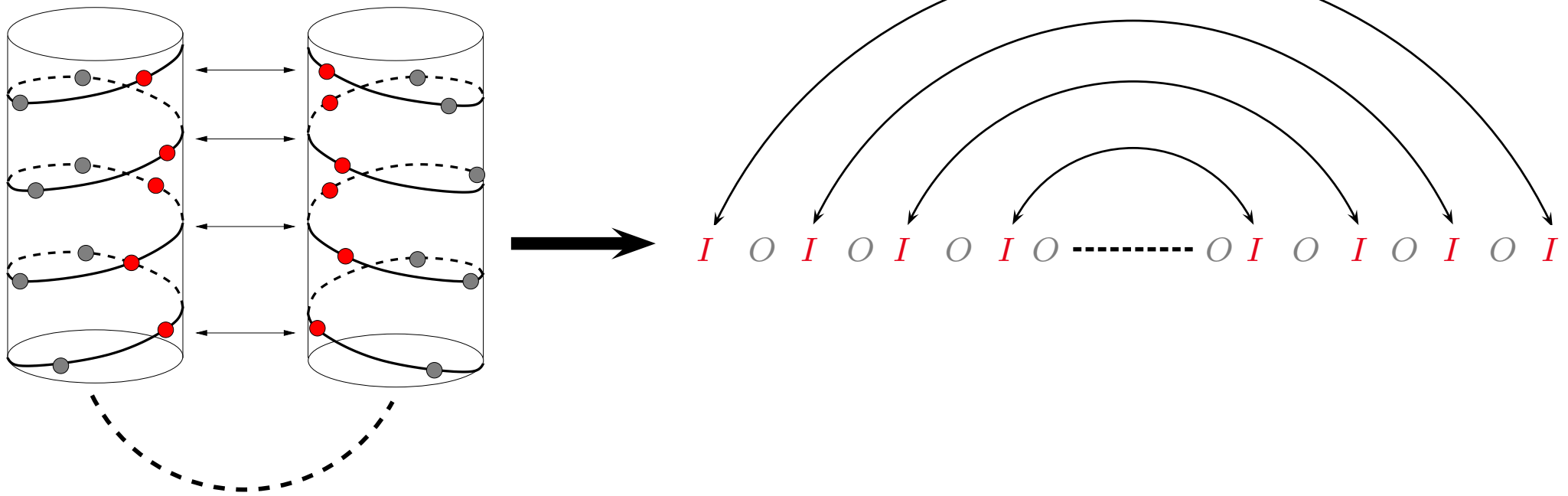
modeling anti-parallel pairing of α -helices



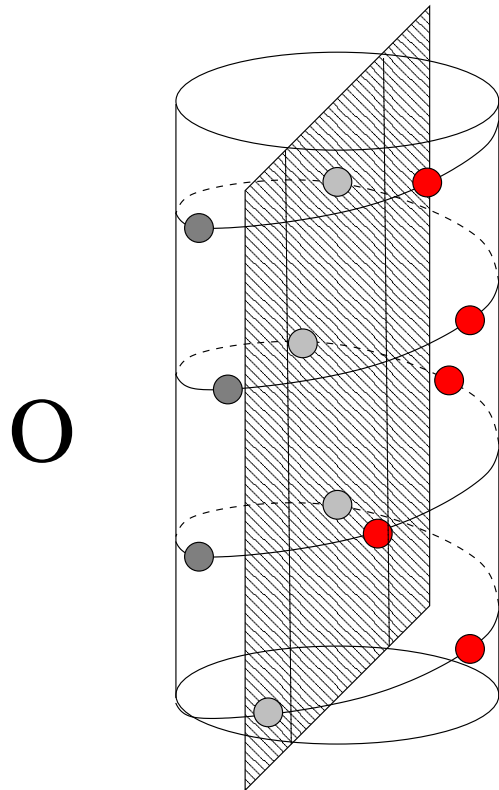
modeling anti-parallel pairing of α -helices



modeling anti-parallel pairing of α -helices



Modeling the local structure of α -helices

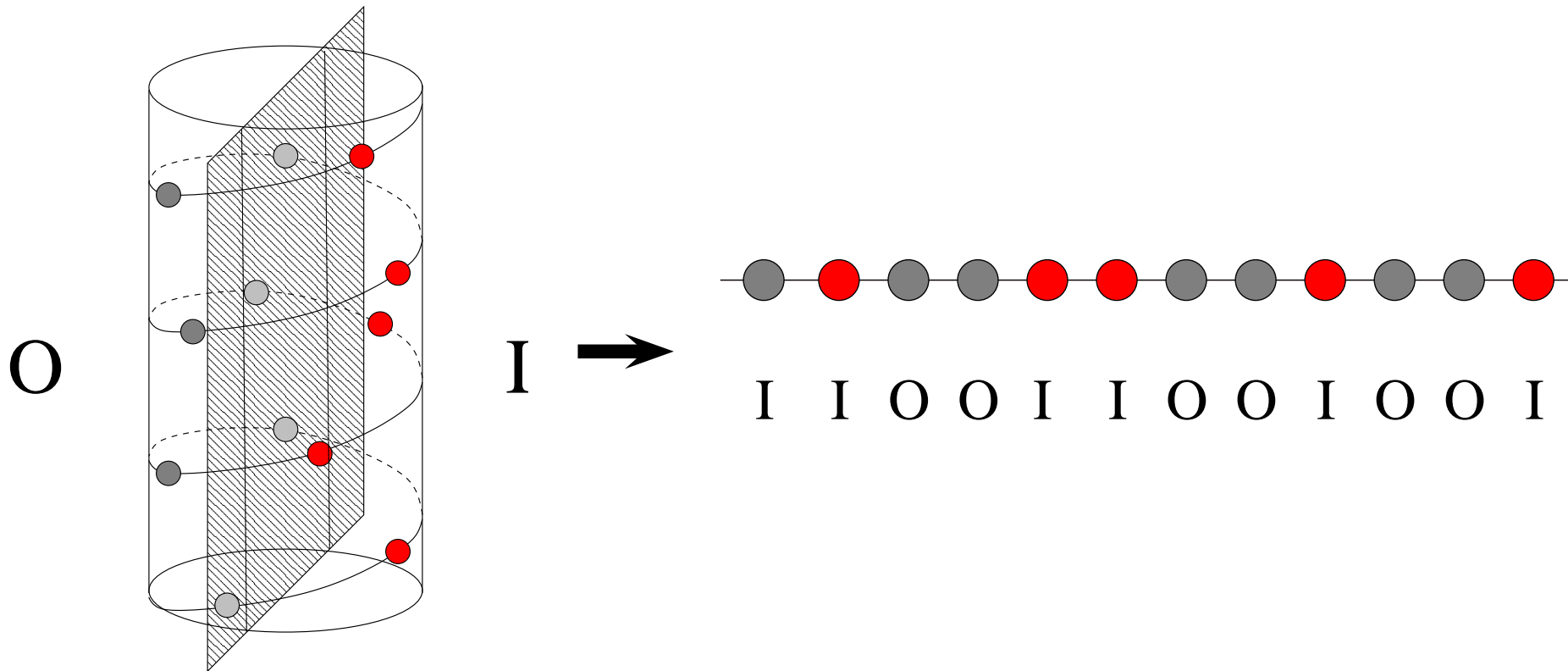


I

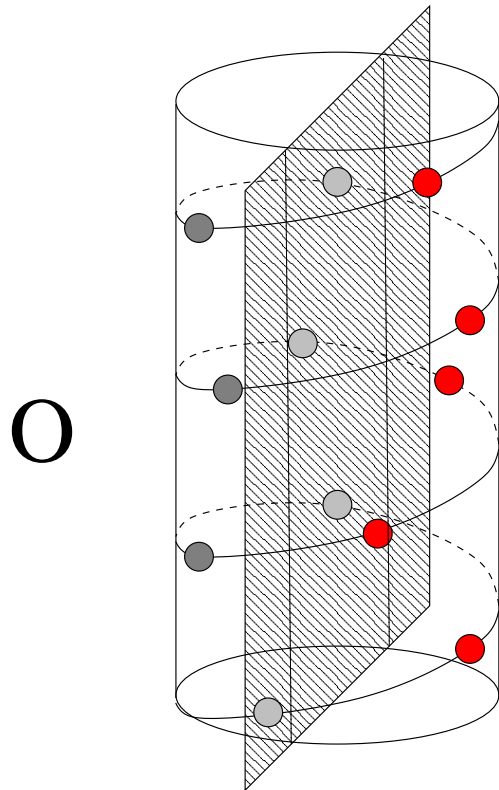


- A *helix* is a stacking of helix turns,
- A *helix turn* is a sequence of residues corresponding to a complete turn around the helix axis,
- A *helical face* is a subsequence of consecutive amino acids of a helix turn,
 - face *I* is involved in pairing,
 - face *O* is opposed to pairing.

Modeling the local structure of α -helices



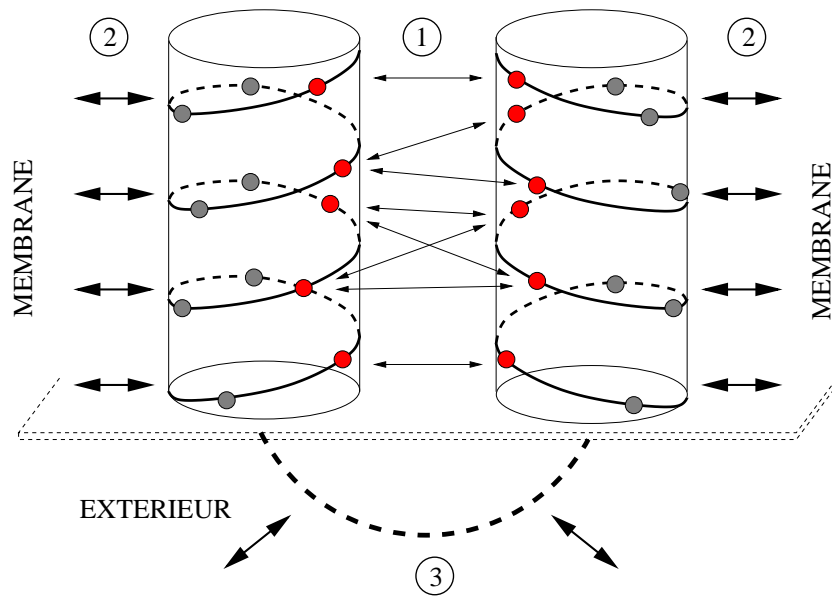
Modeling the local structure of α -helices



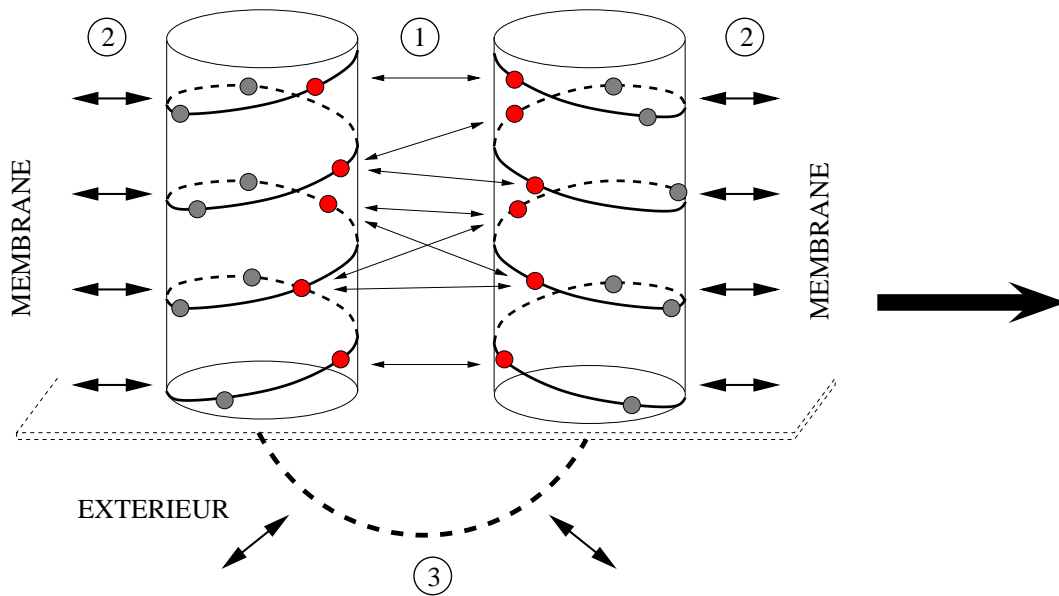
I →

- A helix is an alternate sequence of residues belonging to a *I* or *O* face,
- A helix turn is composed by 3 or 4 consecutive amino acids,
- A helical face has 1 or 2 amino acids,
- On average, 3.6 residus per turn.

Pseudo folding energy



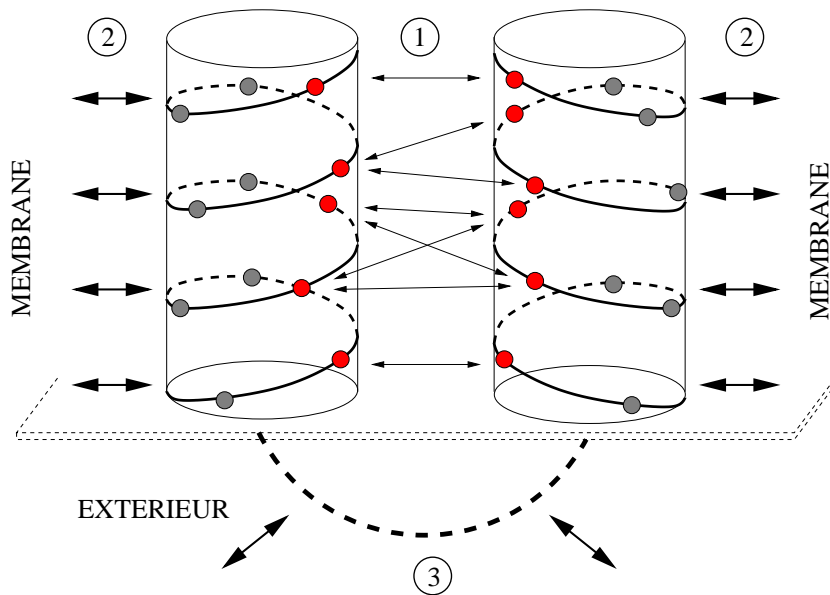
Pseudo folding energy



1. $E_{contact}$ residu interaction energy,

$$E_{contact} = \sum_{i=0}^n f(I_i^k, I_i^{k+1}), \text{ où } f(I_i^k, I_i^{k+1}) = \sum_{\omega_j \in I_i^k} \cdot \sum_{\omega_{j'} \in I_i^{k+1}} \frac{\lambda_j \cdot \lambda_{j'}}{\sqrt{\#I_i^k \cdot \#I_i^{k+1}}}$$

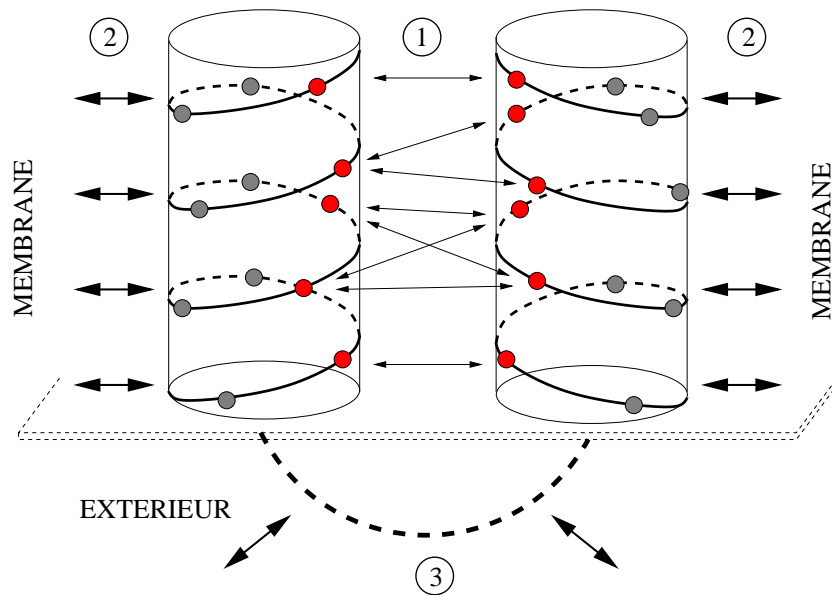
Pseudo folding energy



1. $E_{contact}$ residu interaction energy,
2. E_{memb} membrane interaction energy,

$$E_{memb} = \sum_{\omega_i \in O_i^k} \mathcal{K}_{memb} \cdot \lambda_i$$

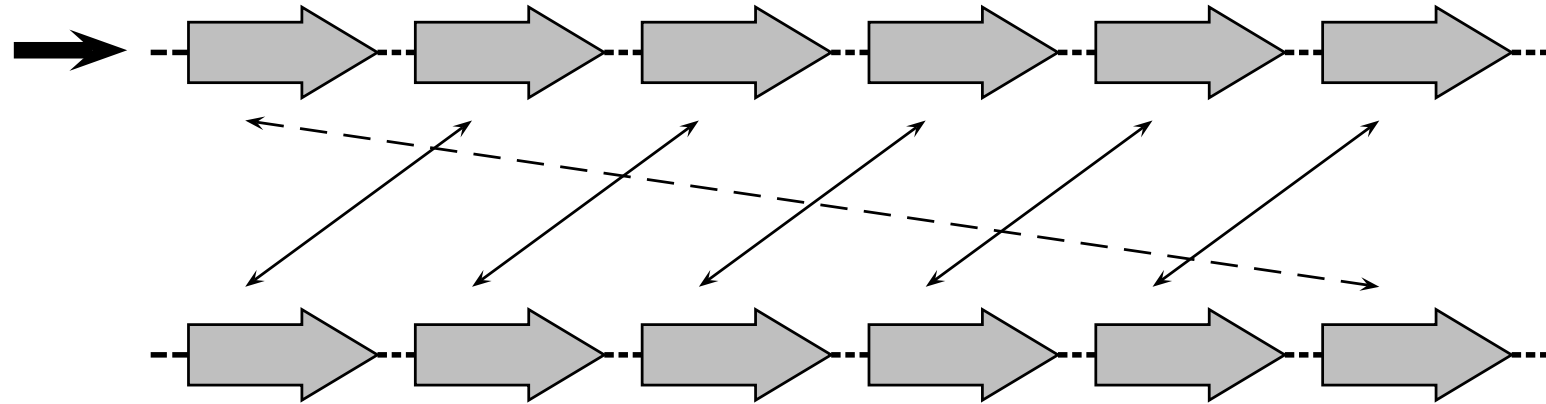
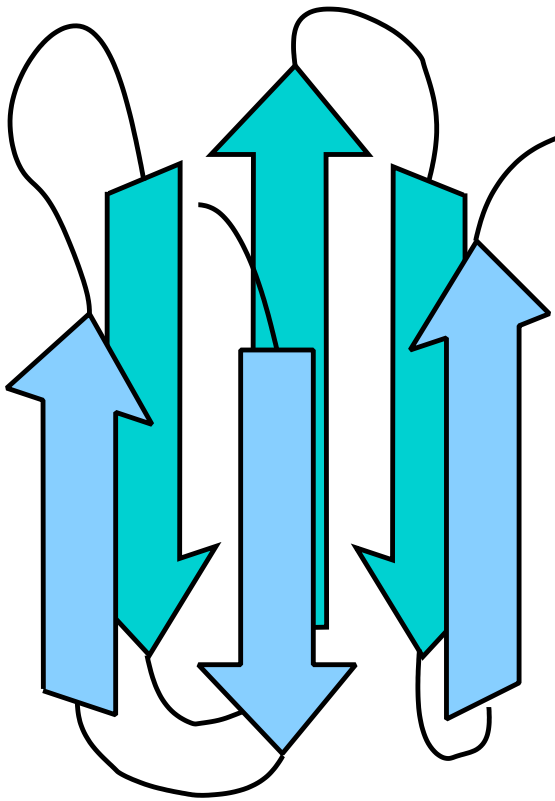
Pseudo folding energy



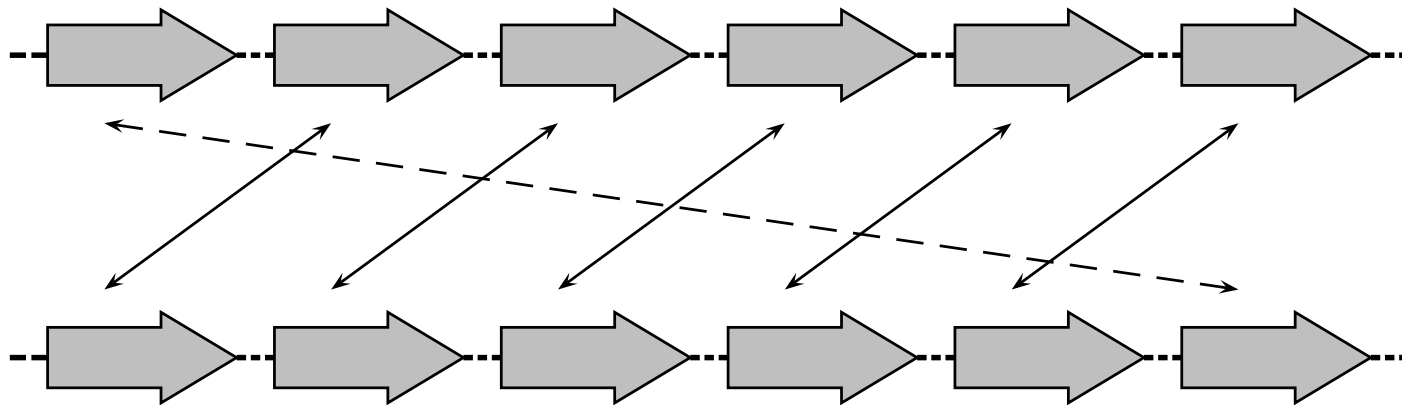
1. $E_{contact}$ residu interaction energy,
2. E_{membr} membrane interaction energy,
3. E_{turn} turn energy.

$$E_{turn} = \mathcal{T}(n - m) + \sum_{i=m}^n \mathcal{K}_{cyt/per} \cdot \lambda_i$$

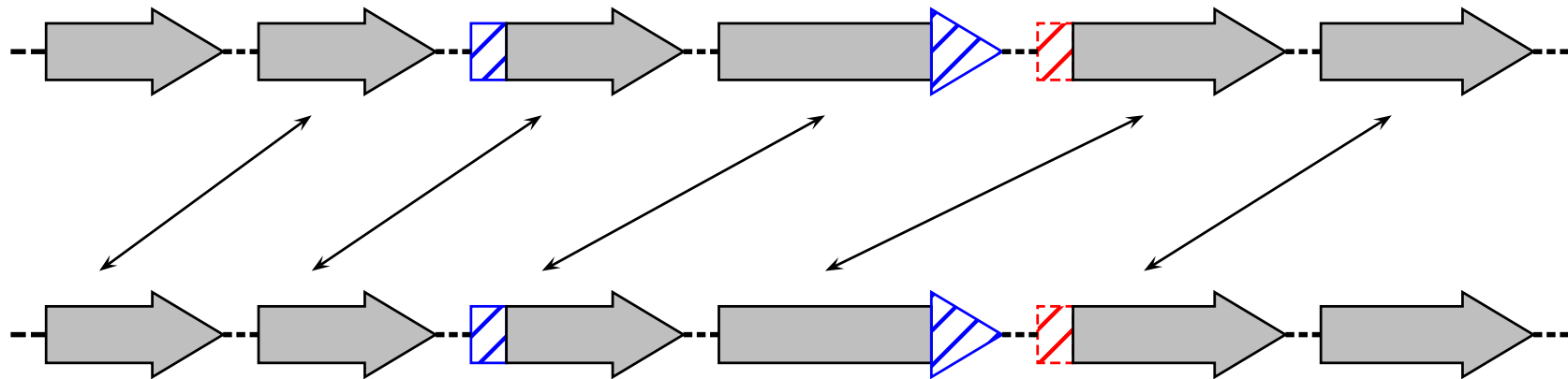
Modeling the overall structure of α -channel



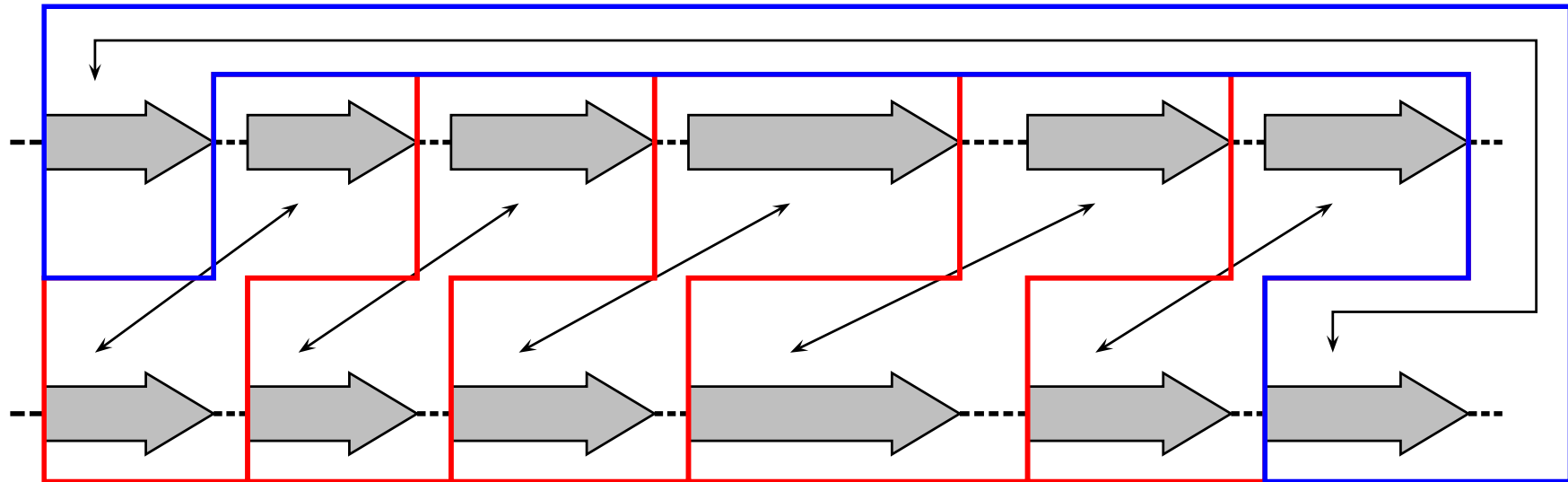
Modeling the overall structure of α -channel



Modeling the overall structure of α -channel



Modeling the overall structure of α -channel



A β -channel is a concatenation of anti-parallel pairings of β -strands.

Grammatical modeling of Transmembrane channels

- Local structure of the secondary structures : rational grammar
- Secondary structure pairing : context-free grammar
- Overall structure of a TM channel : multi-tape context-free grammar
- pseudo folding energy : attributes

Grammatical modeling of TM α -channels

- Regular grammar for α -helix,
- Context-free grammar for α -helix pairings,
- Multi-tape context-free grammar for α -channel,
- Multi-tape S-attribute grammar for α -channel,

Regular grammar for α -helix

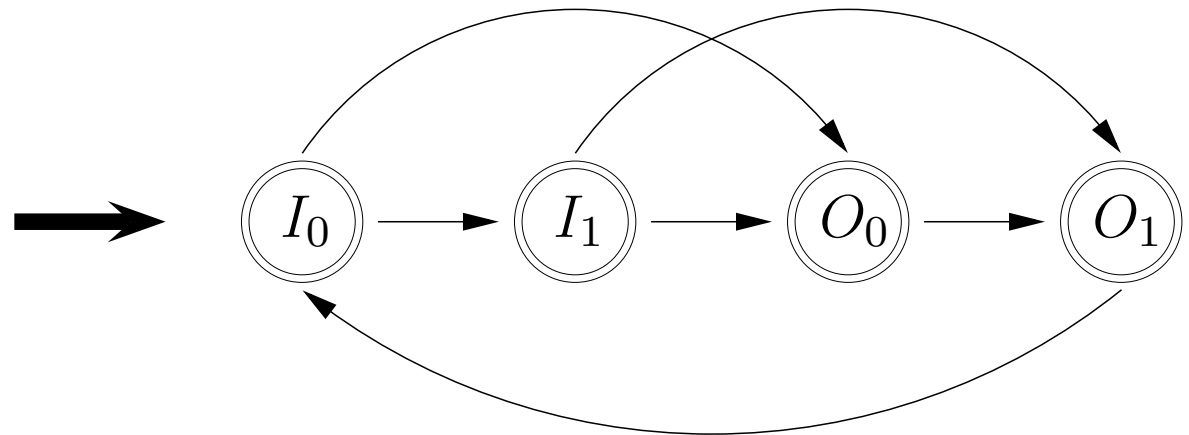
- A helix is an alternate sequence of residues belonging to a I or O face,
- A helix turn is composed by 3 or 4 consecutive amino acids,
- A helical face has 1 or 2 amino acids,
- On average, 3.6 residues per turn.



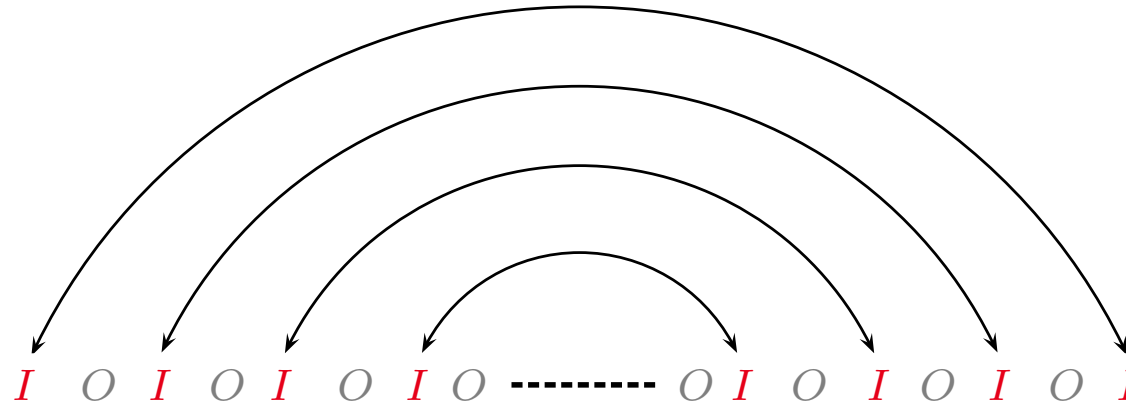
$$P_{helice} = \left\{ \begin{array}{l} S_{helice} \rightarrow I_0 \mid I_1 \mid O_0 \mid O_1 \\ I_0 \rightarrow \bullet I_1 \mid \bullet O_0 \\ I_1 \rightarrow \bullet O_0 \mid \bullet O_1 \\ O_0 \rightarrow \bullet O_1 \\ O_1 \rightarrow \bullet I_0 \end{array} \right.$$

Regular grammar for α -helix

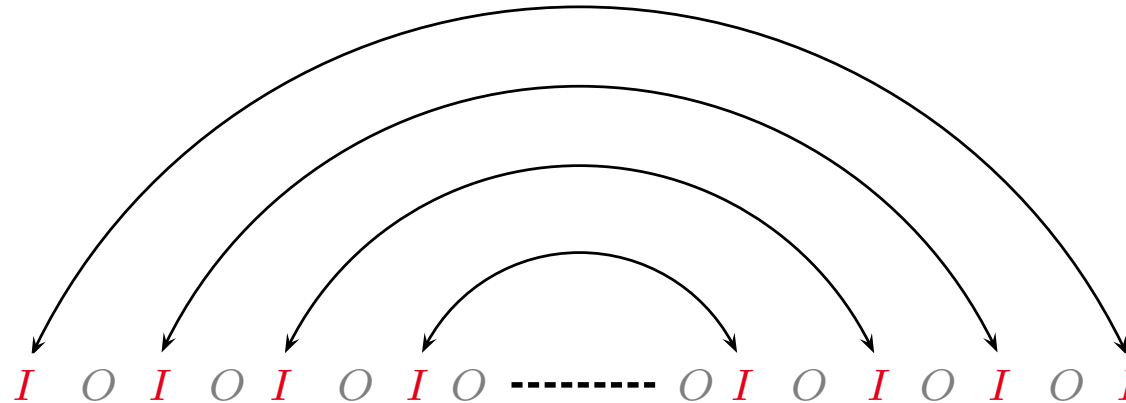
- A helix is an alternate sequence of residues belonging to a I or O face,
- A helix turn is composed by 3 or 4 consecutive amino acids,
- A helical face has 1 or 2 amino acids,
- On average, 3.6 residues per turn.



CFG for α -helix anti-parallel pairings

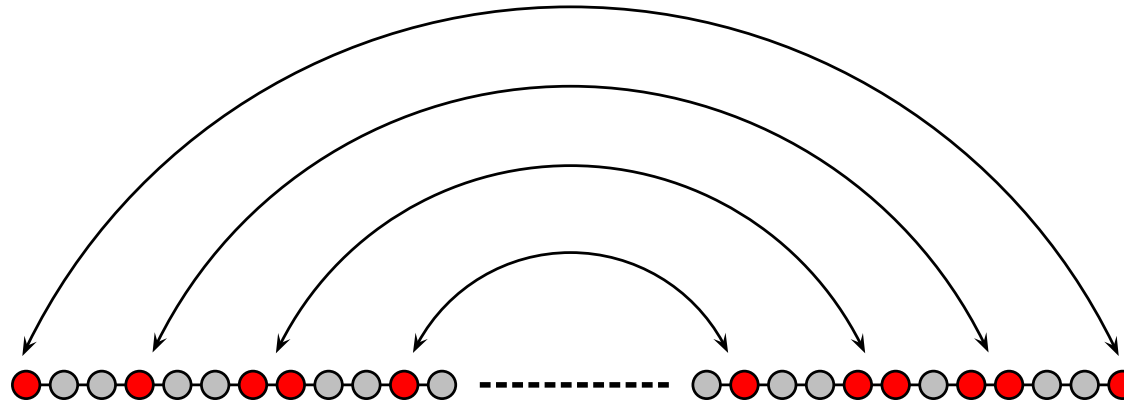


CFG for α -helix anti-parallel pairings



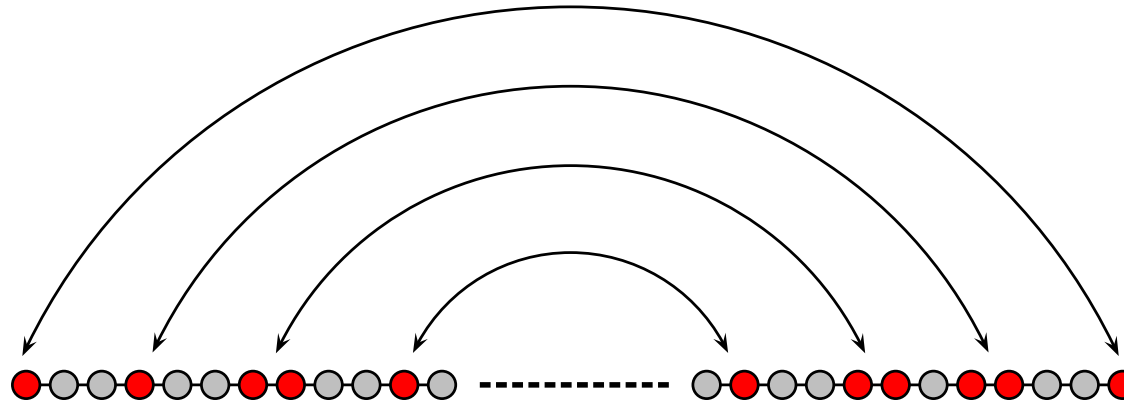
$$P_{ap} = \begin{cases} S_{ap}^{\alpha} & \rightarrow F_O^{\alpha} \mid F_I^{\alpha} \\ F_I^{\alpha} & \rightarrow I F_O^{\alpha} I \mid C_{cyt}^{\alpha} \mid C_{per}^{\alpha} \\ F_O^{\alpha} & \rightarrow O F_I^{\alpha} O \mid C_{cyt}^{\alpha} \mid C_{per}^{\alpha} \\ C_{cyt}^{\alpha} & \rightarrow i C_{cyt}^{\alpha} \mid i \\ C_{per}^{\alpha} & \rightarrow o C_{per}^{\alpha} \mid o \end{cases}$$

CFG for α -helix anti-parallel pairings



$$P_{ap} = \left\{ \begin{array}{l} S_{ap}^{\alpha} \rightarrow F_O^{\alpha} \mid F_I^{\alpha} \\ F_I^{\alpha} \rightarrow I F_O^{\alpha} I \mid C_{cyt}^{\alpha} \mid C_{per}^{\alpha} \\ F_O^{\alpha} \rightarrow O F_I^{\alpha} O \mid C_{cyt}^{\alpha} \mid C_{per}^{\alpha} \\ C_{cyt}^{\alpha} \rightarrow i C_{cyt}^{\alpha} \mid i \\ C_{per}^{\alpha} \rightarrow o C_{per}^{\alpha} \mid o \end{array} \right.$$

CFG for α -helix anti-parallel pairings

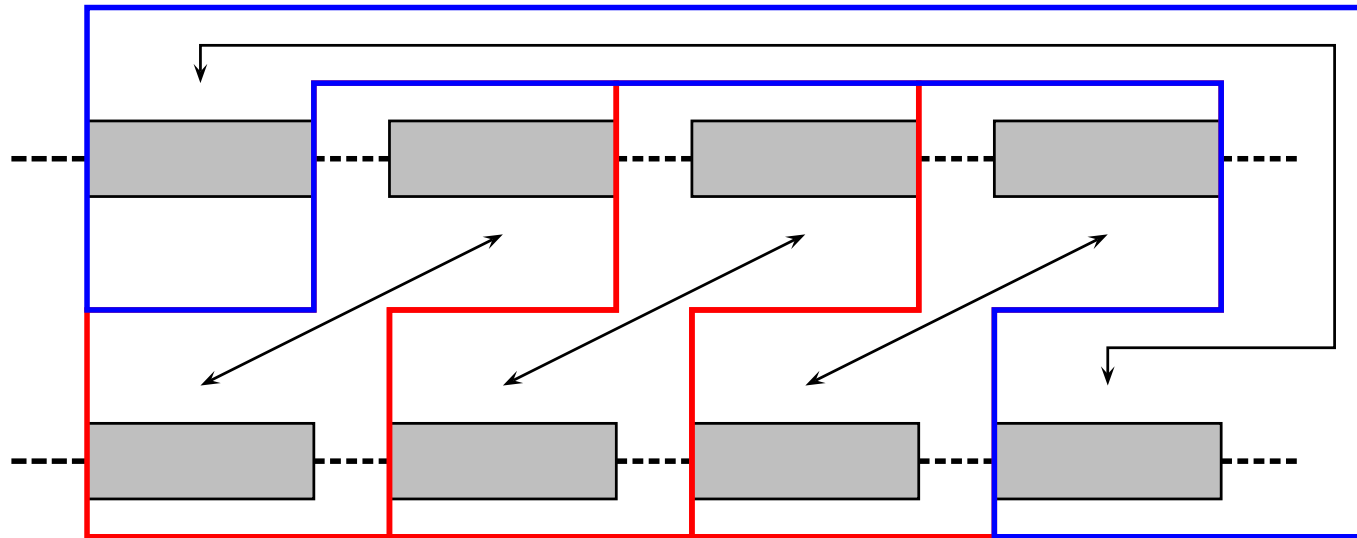


$$P_{ap} = \left\{ \begin{array}{l} S_{ap}^{\alpha} \rightarrow F_O^{\alpha} \mid F_I^{\alpha} \\ F_I^{\alpha} \rightarrow I F_O^{\alpha} I \mid C_{cyt}^{\alpha} \mid C_{per}^{\alpha} \\ F_O^{\alpha} \rightarrow O F_I^{\alpha} O \mid C_{cyt}^{\alpha} \mid C_{per}^{\alpha} \\ C_{cyt}^{\alpha} \rightarrow i C_{cyt}^{\alpha} \mid i \\ C_{per}^{\alpha} \rightarrow o C_{per}^{\alpha} \mid o \end{array} \right. \quad \circ \quad P_{helice} = \left\{ \begin{array}{l} S_{helice} \rightarrow I_0 \mid I_1 \mid O_0 \mid O_1 \\ I_0 \rightarrow \bullet I_1 \mid \bullet O_0 \\ I_1 \rightarrow \bullet O_0 \mid \bullet O_1 \\ O_0 \rightarrow \bullet O_1 \\ O_1 \rightarrow \bullet I_0 \end{array} \right.$$

CFG for α -helix anti-parallel pairings

$$P_{ap}^{\alpha} = \left\{ \begin{array}{l} S_{ap} \rightarrow F_I^{1,1} \mid F_O^{1,1} \qquad 1 \\ F_I^{1,1} \rightarrow \bullet \bullet F_O^{1,1} \bullet \bullet \mid \bullet F_O^{2,1} \bullet \bullet \mid \bullet \bullet F_O^{1,2} \bullet \mid \bullet F_O^{1,1} \bullet \mid C_{cyt}^{\alpha} \qquad 2 \\ F_I^{2,1} \rightarrow \bullet \bullet F_O^{1,1} \bullet \bullet \mid \bullet \bullet F_O^{1,2} \bullet \mid C_{cyt}^{\alpha} \qquad 3 \\ F_I^{1,2} \rightarrow \bullet \bullet F_O^{1,1} \bullet \bullet \mid \bullet F_O^{2,1} \bullet \bullet \mid C_{cyt}^{\alpha} \qquad 4 \\ F_I^{2,2} \rightarrow \bullet \bullet F_O^{2,2} \bullet \bullet \mid C_{cyt}^{\alpha} \qquad 5 \\ F_O^{1,1} \rightarrow \bullet \bullet F_I^{1,1} \bullet \bullet \mid \bullet F_I^{2,1} \bullet \bullet \mid \bullet \bullet F_I^{1,2} \bullet \mid \bullet F_I^{2,2} \bullet \mid C_{cyt}^{\alpha} \qquad 6 \\ F_O^{2,1} \rightarrow \bullet \bullet F_I^{1,1} \bullet \bullet \mid \bullet \bullet F_I^{1,2} \bullet \mid C_{cyt}^{\alpha} \qquad 7 \\ F_O^{1,2} \rightarrow \bullet \bullet F_I^{1,1} \bullet \bullet \mid \bullet F_I^{2,1} \bullet \bullet \mid C_{cyt}^{\alpha} \qquad 8 \\ F_O^{2,2} \rightarrow \bullet \bullet F_I^{1,1} \bullet \bullet \mid C_{cyt}^{\alpha} \qquad 9 \\ C_{cyt}^{\alpha} \rightarrow \bullet C_{cyt}^{\alpha} \mid \bullet \qquad 10 \end{array} \right.$$

MTCFG for α -channels

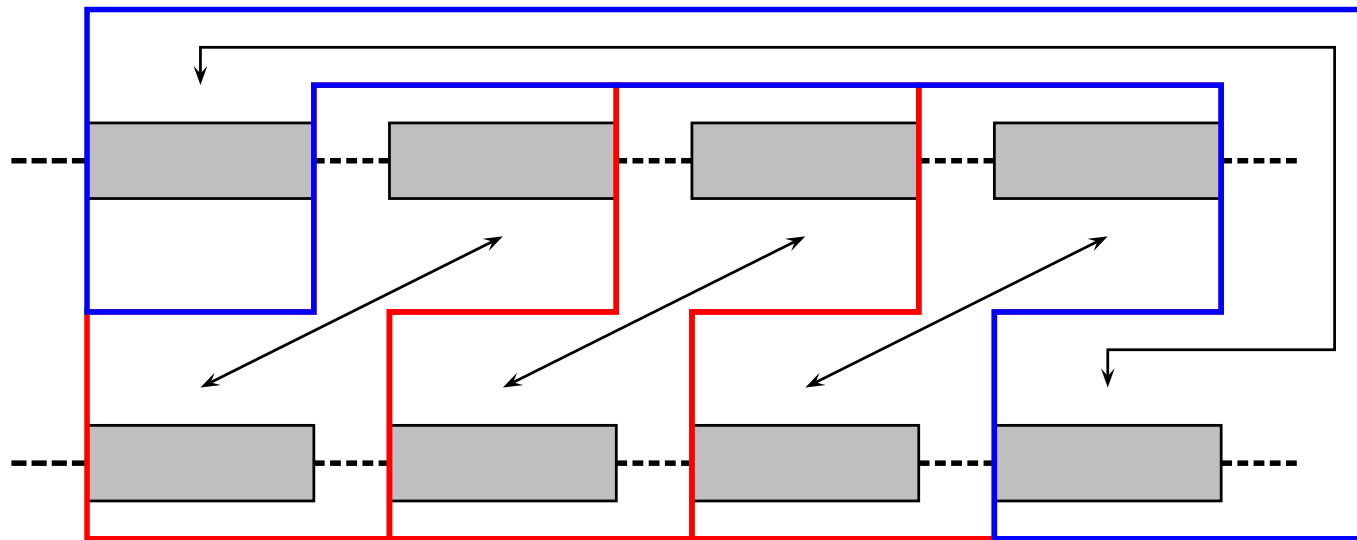


A TM-channel is represented by a 2-tape word :

```

))))))iii)))))oo)))))ii)))))ooo)))))
((((((iii((((((oo((((((ii((((((ooo((((((
    
```

MTCFG for α -channels



A TM-channel is represented by a 2-tape word :

```

))))))-----iii)))))-----oo)))))-----ii)))))-----
-----((((((iii-----((((((oo-----((((((ii-----((((((
    
```


MTCFG for α -channels

$$P_{canal} = \left\{ \begin{array}{llll} S_\alpha & \rightarrow & \begin{bmatrix} t \\ - \end{bmatrix} S_\alpha \begin{bmatrix} - \\ t \end{bmatrix} \mid T_{seq,cyt}^\alpha \mid T_{seq,per}^\alpha & 1 \\ T_{seq,cyt}^\alpha & \rightarrow & T_{cyt}^\alpha T_{seq,per}^\alpha \mid T_{cyt}^\alpha & 2 \\ T_{seq,per}^\alpha & \rightarrow & T_{per}^\alpha T_{seq,cyt}^\alpha \mid T_{per}^\alpha & 3 \\ T_{cyt}^\alpha & \rightarrow & \begin{bmatrix} - \\ t \end{bmatrix} T_{cyt}^\alpha \begin{bmatrix} t \\ - \end{bmatrix} \mid C_{cyt}^\alpha & 4 \\ T_{per}^\alpha & \rightarrow & \begin{bmatrix} - \\ t \end{bmatrix} T_{per}^\alpha \begin{bmatrix} t \\ - \end{bmatrix} \mid C_{per}^\alpha & 5 \\ C_{cyt}^\alpha & \rightarrow & \begin{bmatrix} i \\ i \end{bmatrix} C_{cyt}^\alpha \mid \begin{bmatrix} i \\ i \end{bmatrix} & 6 \\ C_{per}^\alpha & \rightarrow & \begin{bmatrix} o \\ o \end{bmatrix} C_{per}^\alpha \mid \begin{bmatrix} o \\ o \end{bmatrix} & 7 \end{array} \right.$$

How to integrate the pairing rules ?

```
))))iii)))))oo)))))ii)))))oo)))))  
((((iii((((oo((((ii((((oo((((
```

How to integrate the pairing rules ?

MPPMMPMMPPMMiiiPPMPPMMPMMPPooMPPMMPMMPMMiiiPPMPPMMPMMPPooMPPMMPMMPPMM
PMMPPMMPMMPPiiiMPPMMPMMPPMMooPMMPPMPPMMPPIiMPPMMPMMPPMMooMPPMMPMMPPMM

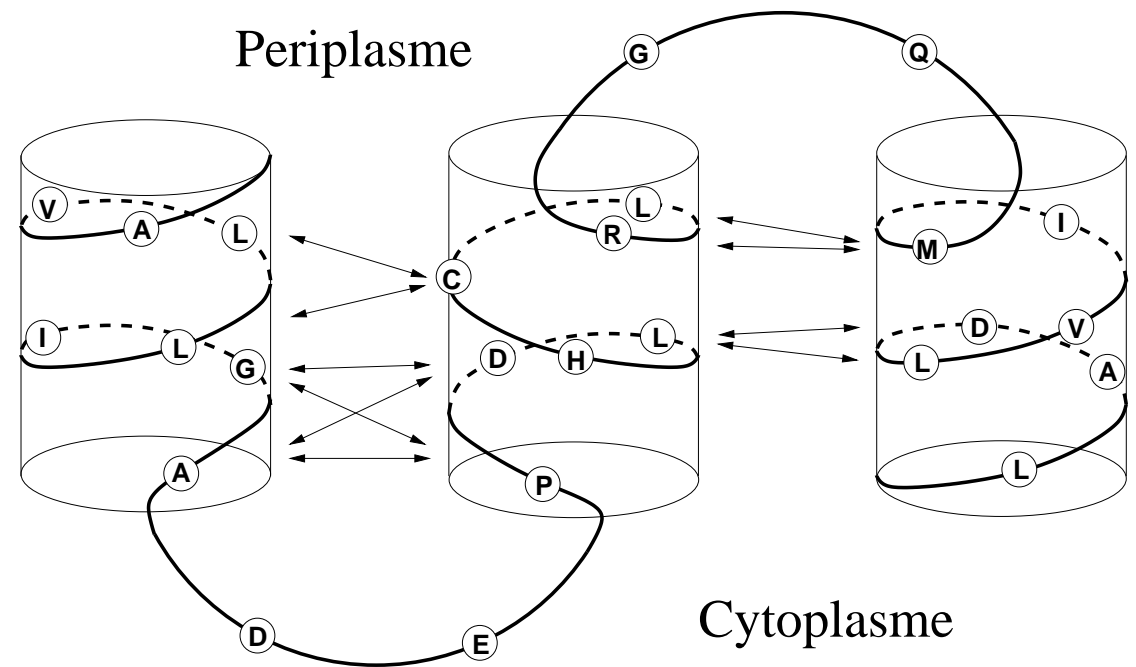
$$P_{ap}^\alpha = \left\{ \begin{array}{ll} S_{ap} \rightarrow F_I^{1,1} \mid F_O^{1,1} & 1 \\ F_I^{1,1} \rightarrow \bullet \bullet F_O^{1,1} \bullet \bullet \mid \bullet F_O^{2,1} \bullet \bullet \mid \bullet \bullet F_O^{1,2} \bullet \mid \bullet F_O^{1,1} \bullet \mid C^\alpha & 2 \\ F_I^{2,1} \rightarrow \bullet \bullet F_O^{1,1} \bullet \bullet \mid \bullet \bullet F_O^{1,2} \bullet \mid C^\alpha & 3 \\ F_I^{1,2} \rightarrow \bullet \bullet F_O^{1,1} \bullet \bullet \mid \bullet F_O^{2,1} \bullet \bullet \mid C^\alpha & 4 \\ F_I^{2,2} \rightarrow \bullet \bullet F_O^{2,2} \bullet \bullet \mid C^\alpha & 5 \\ F_O^{1,1} \rightarrow \bullet \bullet F_I^{1,1} \bullet \bullet \mid \bullet F_I^{2,1} \bullet \bullet \mid \bullet \bullet F_I^{1,2} \bullet \mid \bullet F_I^{2,2} \bullet \mid C^\alpha & 6 \\ F_O^{2,1} \rightarrow \bullet \bullet F_I^{1,1} \bullet \bullet \mid \bullet \bullet F_I^{1,2} \bullet \mid C^\alpha & 7 \\ F_O^{1,2} \rightarrow \bullet \bullet F_I^{1,1} \bullet \bullet \mid \bullet F_I^{2,1} \bullet \bullet \mid C^\alpha & 8 \\ F_O^{2,2} \rightarrow \bullet \bullet F_I^{1,1} \bullet \bullet \mid C^\alpha & 9 \\ C^\alpha \rightarrow \bullet C^\alpha \mid \bullet & 10 \end{array} \right.$$

Multi-tape S-attribute grammar for α -channels

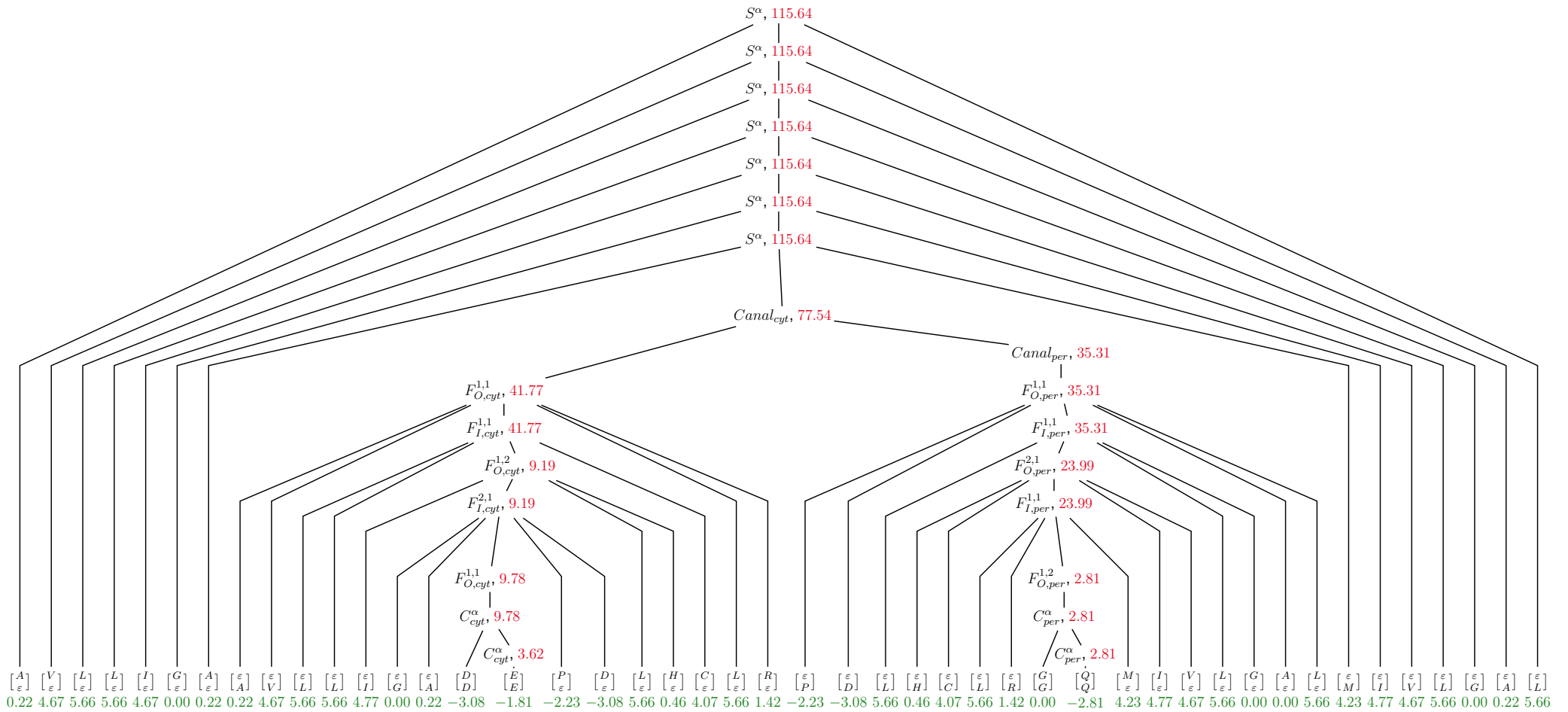
To each production rule, associate a functions which allows a recursive computation of the energy.

An example !

MTSAG for α -channels

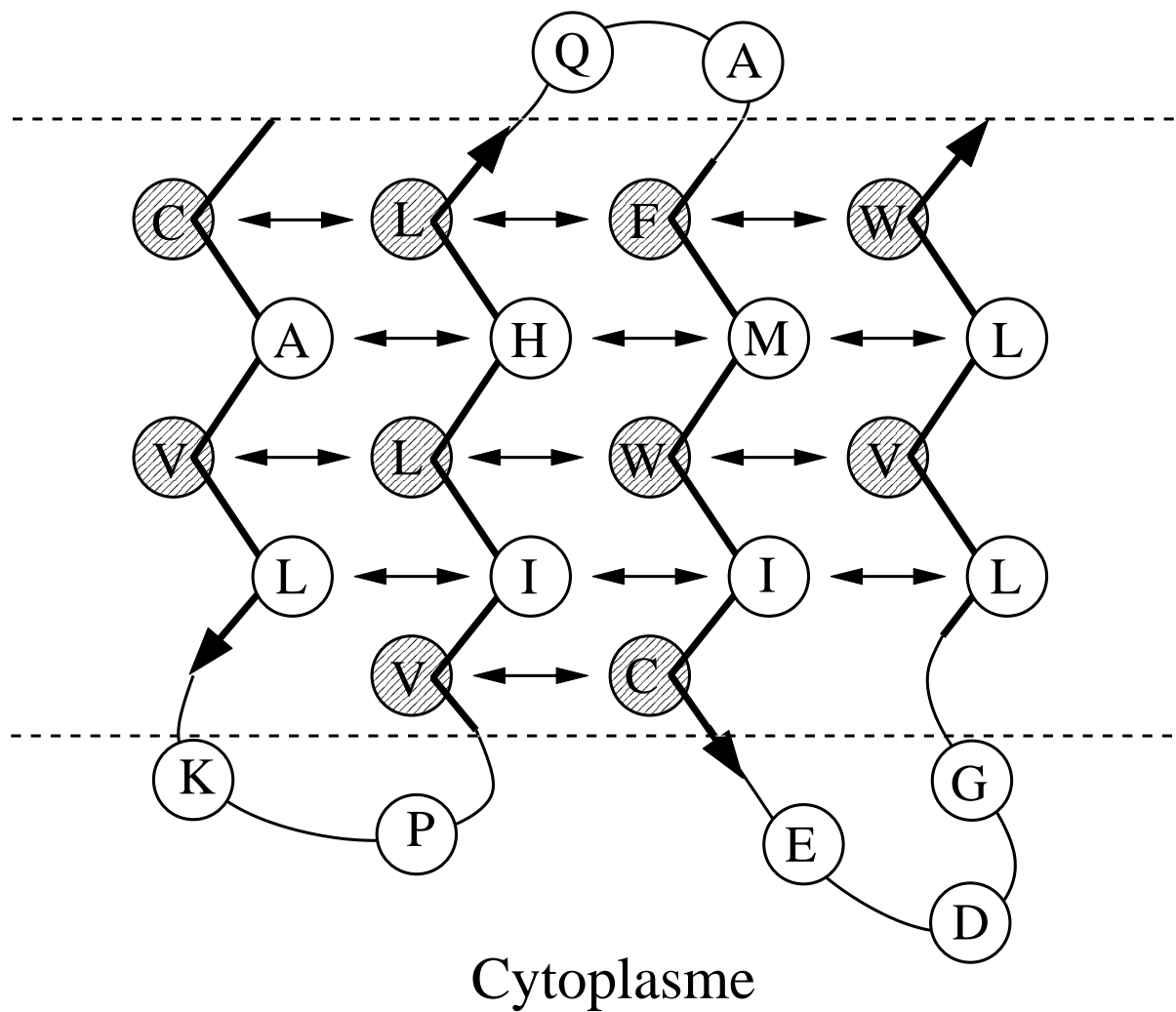


MTSAG for α -channels

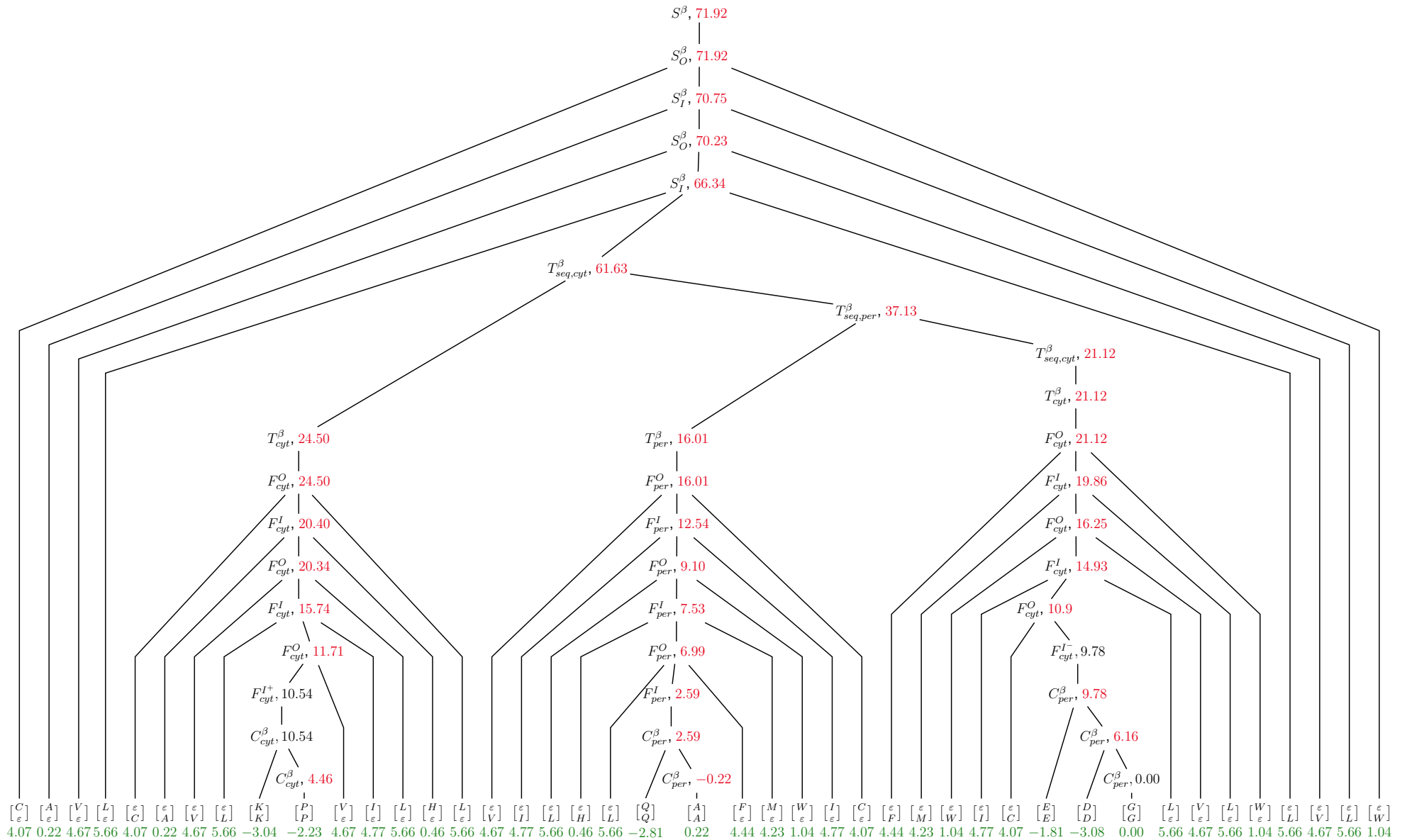


MTSAG for β -channels

Periplasme



MTSAG for β -channels



What has not been said :

- TM-channel closure,
- TM α -helix selection,
- turn selection (between secondary structures),
- constraints on the overlapping of the motifs.

Performance evaluation

- How to realize a structure prediction ?
- How to evaluate a prediction ?
- Results.

How to realize a structure prediction ?

- syntax analysis (GCP algorithm),
- implementation using *mtsag2c* (F. Lefebvre),
- software *tmmtsag...* and now ASTRiD (web interface).

How to realize a structure prediction ?

Example of a β -channel : Porin

```
MAPKDNTWYTGAKLGWSQYHDTGLINNGPTHENKLGAGAFGGYQVNPYVGFEMGYDWLGRMPYKGSVENGAYKAQGVQLTAKLGYPTDDLDIYTRLGG
....TT.EEEEEEEEEES.S.....SS.....EEEEEEEEEE.BTTEEEEEEEEEEE.....SS...EEEEEEEEEEEESSSEEEEEEEEE
.....EEEEEEEEEOOOOOOOOOCBCBCBCBCiIIIIIBCBCBCBCBCCOOOOOOOOOOBCBCBCBCBiIIIIIIIBCBCBCB
.....CBCBCBCBCBOOOOOOOOBCBCBCBCBiIIIIIBCBCBCBCBCCOOOOOOOOOOBCBCBCBCBiIIIIIIIBCBCBCB
```

```
MVWRADTYSNVYGNHDTGVSPVFAGGVEYAITPEIATRLEYQWTNNIGDAHTIGTRPDNGMLSLGVSYRFG
EEEEEEE..SSS..EEEEEEEEEEEEEEEESSSEEEEEEEEEEE.....SS.....EEEEEEEEEE.
CBCBCOOOOOOOOOOOOCBCBCBCBCiIIIBCBCBCBCBCCOOOOOOOOOOOOBCBCBCBCBC.
BCBCBOOOOOOOOOOOOOCBCBCBCBCiIIIBCBCBCBCBCCOOOOOOOOOOOOEEEEEEEEEE.
```

pseudo folding energy : 402.15

How to evaluate a prediction ?

observed :HHHHHHHHHHHHHHHHHHHH..HHHHHHHHHHHHHHHHHH..HHHHHHHHHHHHHHHHHH.....
prediction : ..HHHHHHHHHHHHHHHHHH.....HHHHHHHHHHHH.....HHHHHHHHHHHHHHHHHH..

Definition *A secondary structure is said to be predicted, if it intersects one and only one observed secondary structure.*

Definition *A structure is correctly predicted if all its secondary structures are predicted, almost predicted if the non-predicted secondary structures do not intersect any observed secondary structures, and non-predicted otherwise.*

How to evaluate a prediction ?

.....HHHHHHHHHHHHHHHHHHHH..HHHHHHHHHHHHHHHHHHHH..HHHHHHHHHHHHHHHHHHHH.....
..HHHHHHHHHHHHHHHHHHHH.....HHHHHHHHHHHHHHHHHHHH.....HHHHHHHHHHHHHHHHHHHH.....

How to evaluate a prediction ?

.....HHHHHHHHHHHHHHHHHHHH.....HHHHHHHHHHHHHHHHHHHH.....HHHHHHHHHHHHHHHHHHHH.....
..HHHHHHHHHHHHHHHHHHHH.....HHHHHHHHHHHHHHHHHHHH.....HHHHHHHHHHHHHHHHHHHH.....

Correct

How to evaluate a prediction ?

..... HHHHHHHHHHHHHHHHHHHHH .. HHHHHHHHHHHHHHHHHHHHH .. HHHHHHHHHHHHHHHHHHHHH ..
.. HHHHHHHHHHHHHHHHHHHHH .. HHHHHHHHHHHHHHHHHHHHH .. HHHHHHHHHHHHHHHHHHHHH ..
..... HHHHHHHHHHHHHHHHHHHHH .. HHHHHHHHHHHHHHHHHHHHH .. HHHHHHHHHHHHHHHHHHHHH ..
.. HHHHHHHHHHHHHHHHHHHHH .. HHHHHHHHHHHHHHHHHHHHH ..

Correct

How to evaluate a prediction ?

..... HHHHHHHHHHHHHHHHHHHHHH . . . HHHHHHHHHHHHHHHHHHHHHH . . . HHHHHHHHHHHHHHHHHHHHHH
.. HHHHHHHHHHHHHHHHHHHHHH HHHHHHHHHHHHHHHHHHHHHH HHHHHHHHHHHHHHHHHHHHHH
..... HHHHHHHHHHHHHHHHHHHHHH . . . HHHHHHHHHHHHHHHHHHHHHH . . . HHHHHHHHHHHHHHHHHHHHHH
.. HHHHHHHHHHHHHHHHHHHHHH HHHHHHHHHHHHHHHHHHHHHH HHHHHHHHHHHHHHHHHHHHHH

Correct
Almost

How to evaluate a prediction ?

Estimator for the secondary structure element prediction

$$Q_{ok} = 100 \cdot \frac{\text{number of correctly predicted structures}}{\text{number of proteins}}$$

$$Q_{stm}^{\%obs} = 100 \cdot \frac{\text{number of TM segments correctly predicted}}{\text{number of TM segment observed}}$$

$$Q_{stm}^{\%pred} = 100 \cdot \frac{\text{number of TM segments correctly predicted}}{\text{number of TM segment predicted}}$$

How to evaluate a prediction ?

Estimator for the secondary structure assignment prediction

$$Q_2 = 100 \cdot \frac{\text{number of correctly predicted residus}}{\text{number of residus}}$$

$$Q_{2T}^{\%obs} = 100 \cdot \frac{\text{number of correctly predicted residus in TM segment}}{\text{number of residus observed in TM segments}}$$

$$Q_{2T}^{\%pred} = 100 \cdot \frac{\text{number of correctly predicted residus in TM segment}}{\text{number of residus predicted in TM segments}}$$

$$Q_{2N}^{\%obs} = 100 \cdot \frac{\text{number of correctly predicted residus in non-TM segment}}{\text{number of residus observed in non-TM segments}}$$

$$Q_{2N}^{\%pred} = 100 \cdot \frac{\text{number of correctly predicted residus in non-TM segment}}{\text{number of residus predicted in non-TM segments}}$$

Procedure

- 28 α -TM proteins known at high resolution level,
- 82 α -TM proteins known at low resolution level,
- 14 β -TM proteins known at high resolution level,
- 567 globular proteins,
- computation of sub-optimal structures,
- comparison with 8 other software.

α -TM proteins known at high resolution level

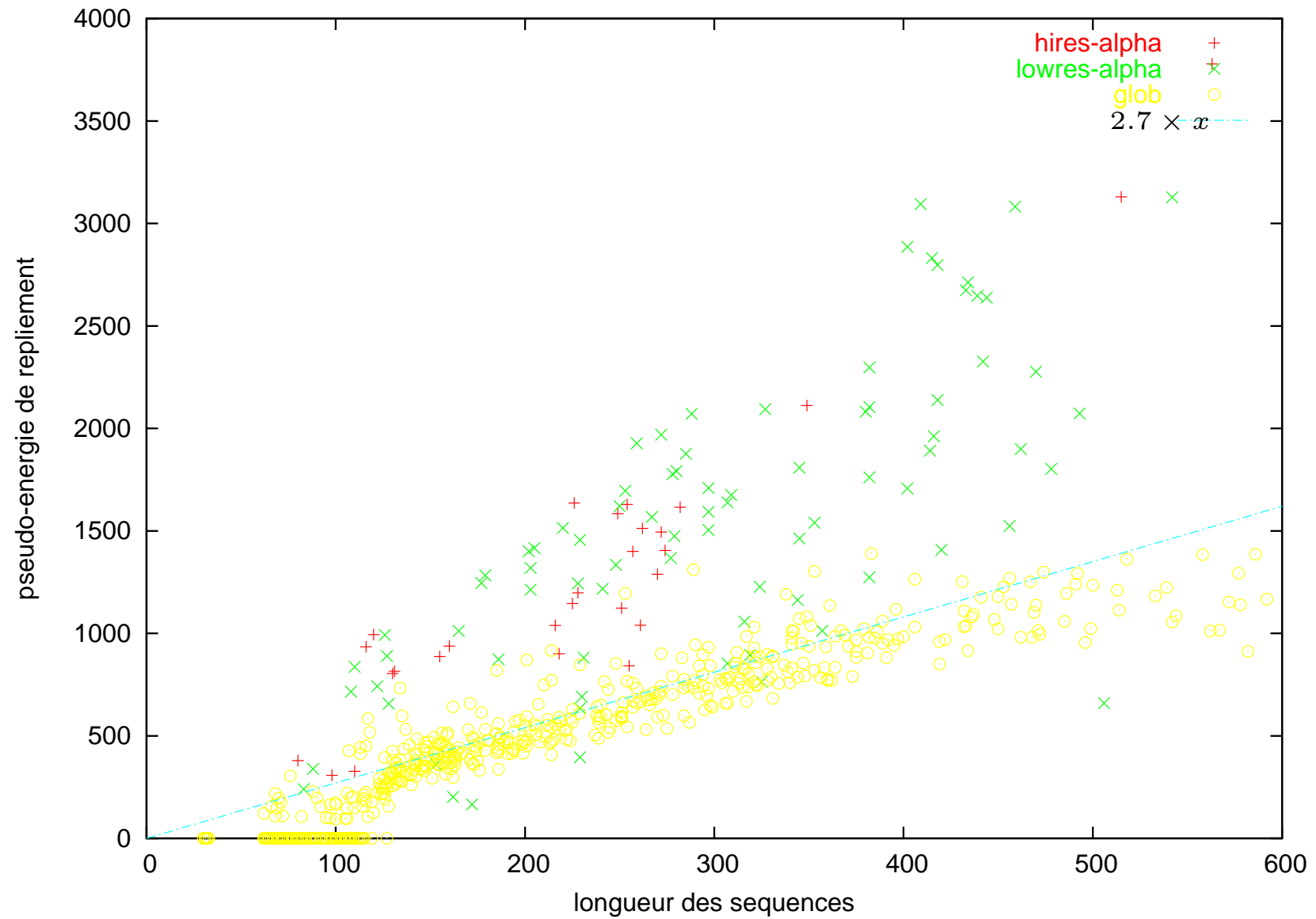
Method	topology	helices		2-states	TM residus		non-TM residus	
	Q_{ok}	$Q_{stm}^{\%obs}$	$Q_{stm}^{\%pred}$	Q_2	$Q_{2T}^{\%obs}$	$Q_{2T}^{\%pred}$	$Q_{2N}^{\%obs}$	$Q_{2N}^{\%pred}$
tmmtsag-basic	75.00(92.86)	97.18	93.88	78.24	87.56	78.77	64.04	77.16
tmmtsag-opt	92.86(100.00)	99.30	98.60	80.06	89.17	80.07	66.17	80.04
hmmtop2	71.43(100.00)	97.18	96.50	79.50	71.84	92.55	91.18	68.00
memsat	60.00(100.00)	94.00	97.92	77.43	67.63	94.34	93.39	63.91
phd-psihtm	71.43(96.43)	88.73	96.92	76.29	68.76	89.55	87.77	64.83
pred-tmr	53.57(100.00)	88.73	100.00	74.41	59.12	97.52	97.71	61.07
sosui	82.14(100.00)	97.18	99.28	80.70	72.71	93.95	92.87	69.10
tmhmm1	71.43(100.00)	96.48	97.16	79.80	72.76	92.13	90.52	68.56
toppred2	71.43(100.00)	95.07	94.41	75.69	66.31	90.98	89.98	63.67

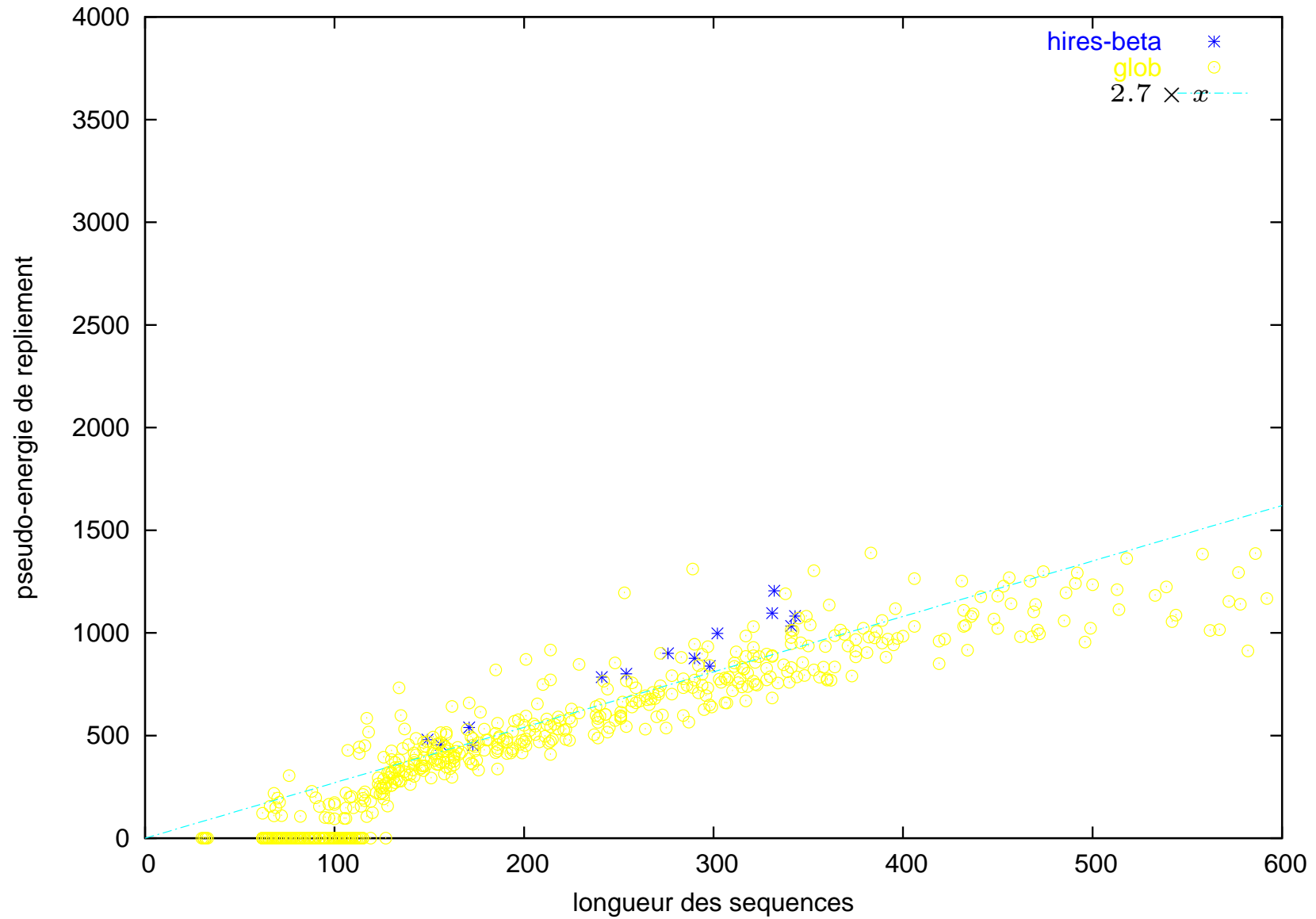
α -TM proteins known at low resolution level

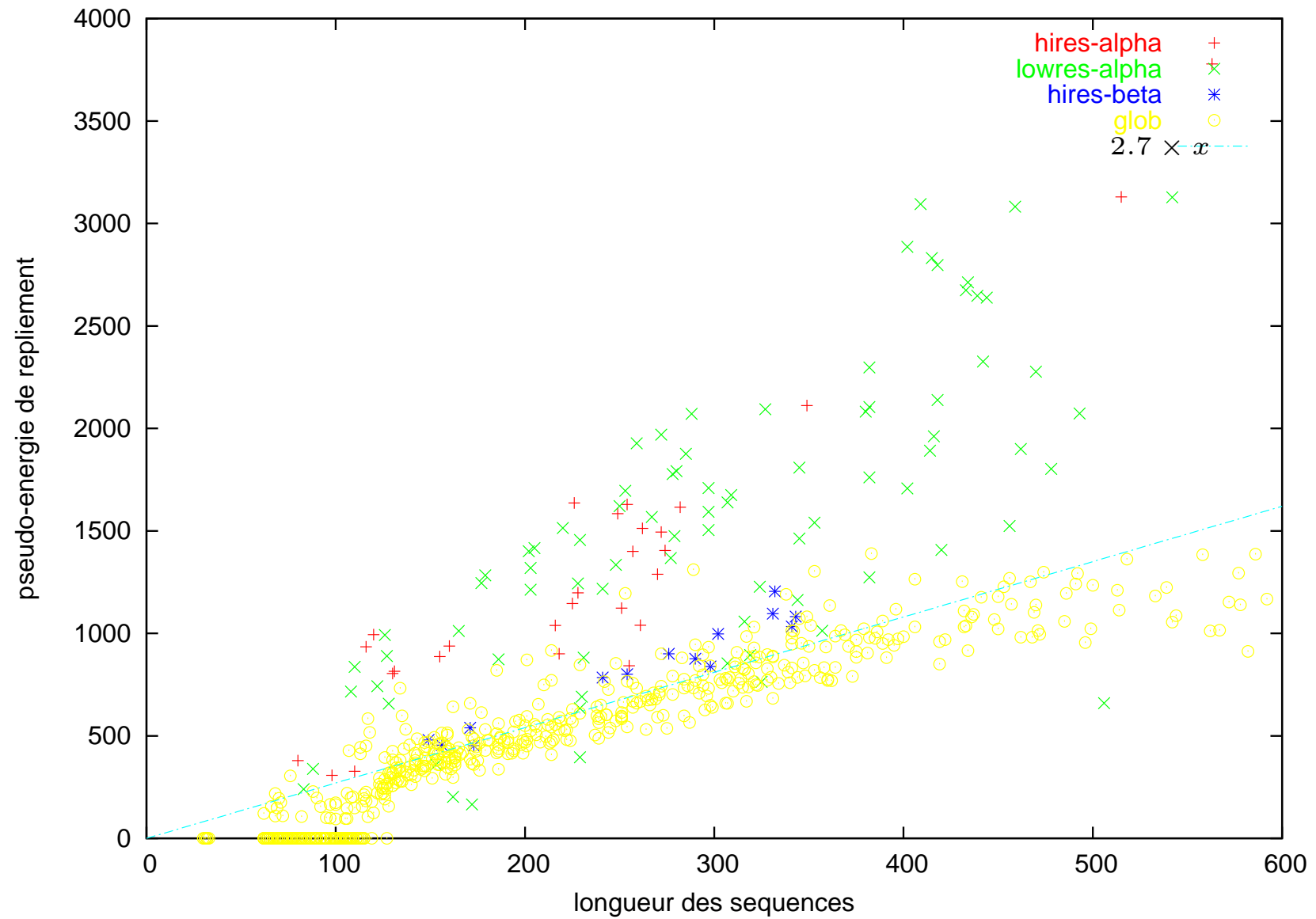
Method	topology	helices		2-states	TM residus		non-TM residus	
	Q_{ok}	$Q_{stm}^{\%obs}$	$Q_{stm}^{\%pred}$	Q_2	$Q_{2T}^{\%obs}$	$Q_{2T}^{\%pred}$	$Q_{2N}^{\%obs}$	$Q_{2N}^{\%pred}$
tmmtsag-basic	37.80(73.17)	87.11	81.68	69.30	87.06	60.31	55.50	84.67
tmmtsag-opt	60.98(95.12)	95.70	91.42	75.25	90.42	65.79	63.47	89.51
hmmtop2	60.98(86.59)	87.11	93.31	84.39	81.69	82.44	86.48	85.88
memsat	54.88(92.68)	91.02	93.20	85.26	81.40	84.80	88.35	85.60
phd-psihtm	24.36(43.59)	60.04	72.02	83.99	91.08	76.60	78.51	91.93
pred-tmr	50.00(96.34)	88.09	96.57	85.19	75.86	88.63	92.44	83.14
sosui	48.78(90.24)	86.33	94.85	82.36	79.17	80.21	84.83	83.98
tmhmm1	69.51(89.02)	90.23	95.45	85.32	83.02	83.34	87.11	86.85
toppred2	53.66(86.59)	83.59	95.75	83.46	74.48	85.81	90.43	82.02

β -TM proteins known at high resolution level

Method	topology	strands		2-states	TM residus		non-TM residus	
	Q_{ok}	$Q_{stm}^{\%obs}$	$Q_{stm}^{\%pred}$	Q_2	$Q_{2T}^{\%obs}$	$Q_{2T}^{\%pred}$	$Q_{2N}^{\%obs}$	$Q_{2N}^{\%pred}$
total								
tmmtsag-basic	7.14(78.57)	83.14	66.82	64.87	72.71	69.99	53.07	56.37
tmmtsag-opt	21.43(92.86)	90.70	82.98	66.80	74.03	71.66	55.93	58.85
tmb-hmm	64.29(100.00)	97.09	97.09	84.04	85.75	87.45	81.47	79.15
small proteins								
tmmtsag-basic	50.00(75.00)	90.62	80.56	70.23	74.30	79.50	62.21	55.10
tmmtsag-opt	50.00(100.00)	93.75	93.75	76.12	78.97	84.08	70.51	62.96
tmb-hmm	75.00(100.00)	96.88	96.88	82.33	83.64	89.05	79.72	71.19







Method	false negatives (%)	false positives (%)			
		$\Delta_{hires}^{tm,\alpha}$	$\Delta_{lowres}^{tm,\alpha}$	$\Delta_{hires}^{tm,\beta}$	$\Delta_{\alpha+\beta}^{tm}$
tmmtsag-basic	14.8	0	8.5	7.1	6.2
hmmtop2	6	0	1	-	-
phd-psihtm	2	3	8	-	-
pred-tmr	4	8	1	-	-
sosui	1	8	4	-	-
tmhmm1	1	8	4	-	-
toppred2	10	8	11	-	-
tmb-hmm	10	-	-	0	-

What has been done :

tmmtsag is the first software able to :

- use efficiently and without restrictions, long range interactions,
- unify α -channel and β -channel in the same model,
- discriminate the 3 categories α , β and globular.

Advantages :

- simple, versatile and efficient,
- no learning method used (more able to detect new structures, with lack of experimental data's).

What has to be done :

- refining the physical approximate model (structure and energy),
- integrating existing methods,
- modeling the quaternary structure,
- extension to globular proteins.

What we are doing :

- model refinement (joint work with T. Simonson),
- study of single point mutations in human rhodopsin (joint work with P. Clote),
- reconstructing the tertiary structure from contact predictions,
- screening of a complete genome,
- web interface : ASTRiD
<http://www.lix.polytechnique.fr/Labo/Jerome.Waldispuhl/astrid/>.

This work :

- J. Waldispühl and J.-M. Steyaert,
Modeling and Predicting All- α Transmembrane Proteins Including
Helix-helix Pairing, *to appear in Theoretical Computer Science, special
issue on pattern discovery in the post genom*, 26 pages, 2005.
available online : <http://www.elsevier.com/locate/tcs/>

Multi-tape S-attributed grammars :

- F. Lefebvre,
A Grammar-Based Unification of Several Alignment and Folding
Algorithms, *Proceedings of the Fourth International Conference on
Intelligent Systems for Molecular Biology*, pp 143-154, 1996.

Slides :

These slides have been realized with the Youpla \LaTeX package, provided by
E. Thome. Download it, at : <http://www.loria.fr/~thome/>