

# Using residue coevolution to retrieve protein homologs<sup>1</sup>

## ComPotts

Hugo Talibart, François Coste



Co-evolutionary methods for the prediction and design  
of protein structure and interactions

CECAM-HQ-EPFL, June 18, 2019

---

<sup>1</sup>work in progress

# The Dyliss bioinformatics team

<http://www.irisa.fr/dyliss>





- Symbiose (bioinformatics Irisa/Inria Rennes):

- Dyliss research team
- Genscale research team
- Genouest bioinformatics platform

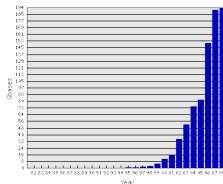
Seminars: <http://symbiose.irisa.fr/symbioseSeminars>

- Biogenouest western France life science and environment network  
Marine biology, agriculture/food-processing, human health, and bioinformatics.

# Motivation

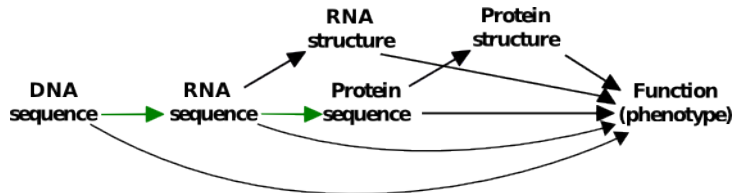
# Sequences annotation problem

High throughput production of raw sequences



## Problem

Function(s) of these sequences ?



# Protein function?

- *In-vivo / in-vitro* experiments

Especially on model organisms:

- Gene knockout and others mutations  
→ key sequence(s) for a function
- Structure determination  
→ key positions for a function
- ...



Does not scale well. . .

- To face the (ever-increasing) amount of available sequences, automatic methods are needed  $\rightsquigarrow$  *in-silico* functional or structural predictions.

Classical approach to predict the function of a new gene sequence

Search for annotated homologs . . .

# Retrieve the homologs of a protein gene

## Search for a significant match with:

- an (already annotated) protein sequence, e.g. with BLAST<sup>2</sup>

>1shg:A

AKELVLALYDYQEKSPREVTMKKGDILTLLNSTNKDWWKVEVNDRQGFVPAAYVKKLD

---

<sup>2</sup>S. F. Altschul et al. “Basic local alignment search tool”. *Journal of molecular biology* (1990).

<sup>3</sup>S. R. Eddy. “Profile hidden Markov models”. *Bioinformatics* 14.9 (1998), pp. 755–763.

<sup>4</sup>M. Steinegger et al. “HH-suite3 for fast remote homology detection and deep protein annotation”. *bioRxiv* (2019), p. 560029.

# Retrieve the homologs of a protein gene



## Search for a significant match with:

- an (already annotated) protein sequence, e.g. with BLAST<sup>2</sup>  
>1shg:A  
AKELVLALYDYQEKSPREVTMKKGDILTLLNSTNKDWWKVEVNDRQGFVPAAYVKKLD

---

<sup>2</sup>S. F. Altschul et al. “Basic local alignment search tool”. *Journal of molecular biology* (1990).

<sup>3</sup>S. R. Eddy. “Profile hidden Markov models”. *Bioinformatics* 14.9 (1998), pp. 755–763.

<sup>4</sup>M. Steinegger et al. “HH-suite3 for fast remote homology detection and deep protein annotation”. *bioRxiv* (2019), p. 560029.

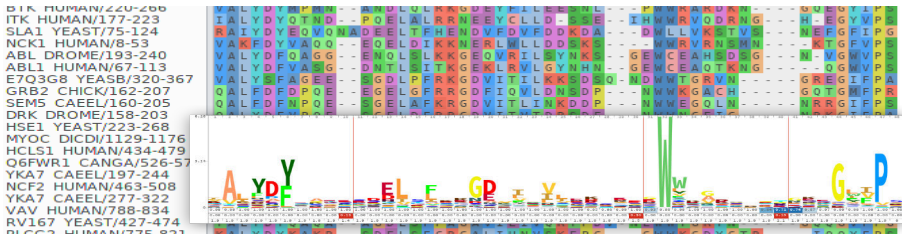


# Retrieve the homologs of a protein gene



## Search for a significant match with:

- an (already annotated) protein sequence, e.g. with BLAST<sup>2</sup>  
>1shg:A  
AKELVLALYDYQEKS PREVTM KKGDI L TLLNSTNKD WVKVEV NDRQGFVPAAYVKKLD
- the Profile HMM of a protein family, e.g. with HMMER<sup>3</sup> or HH-suite<sup>4</sup>



<sup>2</sup>S. F. Altschul et al. “Basic local alignment search tool”. *Journal of molecular biology* (1990).

<sup>3</sup>S. R. Eddy. “Profile hidden Markov models”. *Bioinformatics* 14.9 (1998), pp. 755–763.

<sup>4</sup>M. Steinegger et al. “HH-suite3 for fast remote homology detection and deep protein annotation”. *bioRxiv* (2019), p. 560029.

Update on Predicting the Function of Unknown Proteins

Proteins of Unknown Biochemical Function: A Persistent Problem and a Roadmap to Help Overcome It<sup>1</sup>

Thomas D. Niehaus<sup>2</sup>, Antje M.K. Thamm<sup>2</sup>, Valérie de Crécy-Lagard, and Andrew D. Hanson<sup>1</sup>  
Horticultural Sciences Department (T.D.N., A.M.K.T., A.D.H.) and Microbiology and Cell Science Department (V.d.C.-L.), University of Florida, Gainesville, Florida 32611  
ORCID ID: 0000-0003-2585-4540 (A.D.H.)

RESEARCH ARTICLE

An Approach to Function Annotation for Proteins of Unknown Function (PUFs) in the Transcriptome of Indian Mulberry

K. H. Dhanyalakshmi<sup>1</sup>, Mahantesh B. N. Naika<sup>2\*</sup>, R. S. Sajeevan<sup>1</sup>, Gommen K. Mathew<sup>2</sup>, K. Mohamed Shas<sup>2</sup>, Ramarathnan Sowdhamini<sup>2\*</sup>, Karaba N. Nataraja<sup>1</sup>  
<sup>1</sup> Department of Crop Physiology, University of Agricultural Sciences, GKVK, Bengaluru, 560005, India, <sup>2</sup> National Centre for Biological Sciences, TIFR, GKVK campus, Bengaluru, Karnataka, India

About 16 and 30% of proteins are unannotated in bacteria and yeast genomes. In eukaryotes, over 40% of the proteins encoded by genomes is reported to lack functional annotation

PLOS ONE | DOI:10.1371/journal.pone.0151323 March 16, 2016

Even in Arabidopsis (*Arabidopsis thaliana*), only approximately 40% of enzyme- and transporter-encoding genes have credible functional annotations, and this number is even lower in nonmodel plants

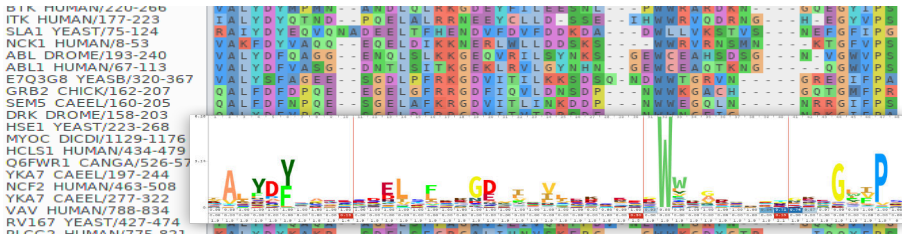
*Plant Physiology*<sup>®</sup>, November 2015, Vol. 169, pp. 1436–1442

# Retrieve the homologs of a protein gene



## Search for a significant match with:

- an (already annotated) protein sequence, e.g. with BLAST<sup>2</sup>  
>1shg:A  
AKELVLALYDYQEKSPEVMTMKKGDILTLLNSTNKDWKVEVNDNRQGFVPAAYVKKLD
- the Profile HMM of a protein family, e.g. with HMMER<sup>3</sup> or HH-suite<sup>4</sup>

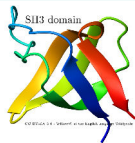


<sup>2</sup>S. F. Altschul et al. “Basic local alignment search tool”. *Journal of molecular biology* (1990).

<sup>3</sup>S. R. Eddy. “Profile hidden Markov models”. *Bioinformatics* 14.9 (1998), pp. 755–763.

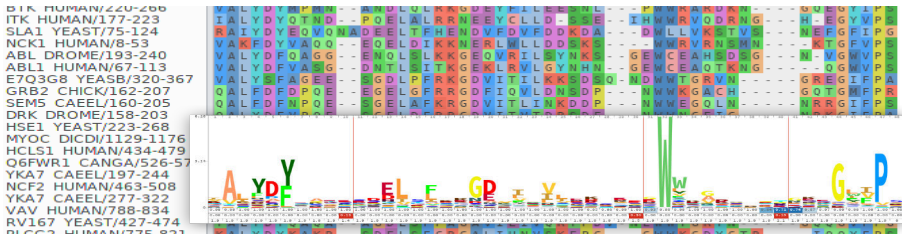
<sup>4</sup>M. Steinegger et al. “HH-suite3 for fast remote homology detection and deep protein annotation”. *bioRxiv* (2019), p. 560029.

# Retrieve the homologs of a protein gene



## Search for a significant match with:

- an (already annotated) protein sequence, e.g. with BLAST<sup>2</sup>  
>1shg:A  
AKELVLALYDYQEKS PREVTM KKGDI LTLN STNKD WVKVEV NDRQGFVPAAYVKKLD
- the Profile HMM of a protein family, e.g. with HMMER<sup>3</sup> or HH-suite<sup>4</sup>



Score each position independently :-)

<sup>2</sup>S. F. Altschul et al. "Basic local alignment search tool". *Journal of molecular biology* (1990).

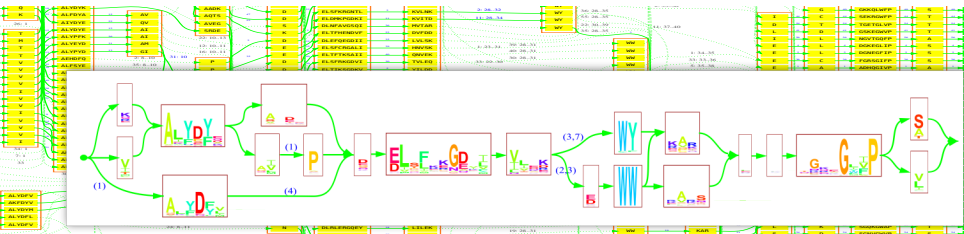
<sup>3</sup>S. R. Eddy. "Profile hidden Markov models.". *Bioinformatics* 14.9 (1998), pp. 755–763.

<sup>4</sup>M. Steinegger et al. "HH-suite3 for fast remote homology detection and deep protein annotation". *bioRxiv* (2019), p. 560029.



## Automatic characterization of protein sequence families with:

- Automata (Protomata-Learner<sup>5,6</sup>)



Local dependencies :-)

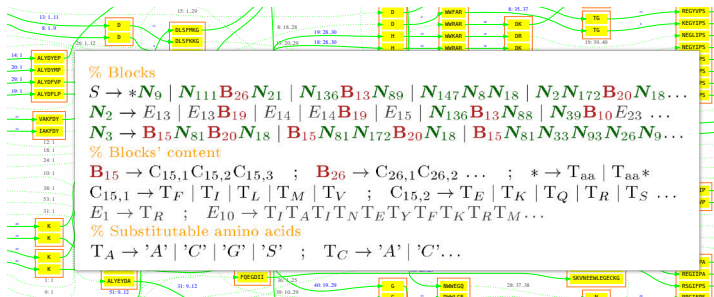
<sup>5</sup>G. Kerbellec. “Apprentissage d’automates modélisant des familles de séquences protéiques”. PhD thesis. Université de Rennes 1, Apr. 2008, p. 139.

<sup>6</sup>A. Bretaudeau et al. “CyanoLyase: a database of phycobilin lyase sequences, motifs and functions”. *Nucleic Acids Research* 41.Database-Issue (2013), pp. 396–401.



## Automatic characterization of protein sequence families with:

- Context-free grammars (ReGLiS<sup>7</sup>, see also<sup>8</sup>)

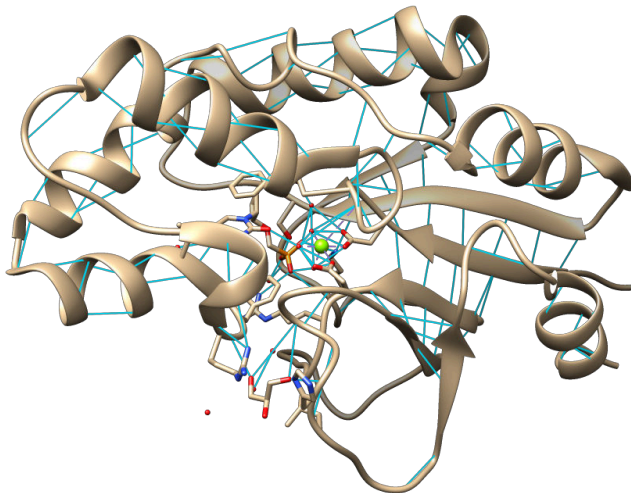


Nested dependencies :-D

<sup>7</sup>F. Coste, G. Garet, and J. Nicolas. "A bottom-up efficient algorithm learning substitutable languages from positive examples". *ICGI*. 2014.

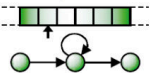
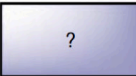
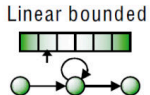


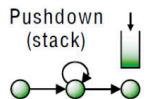




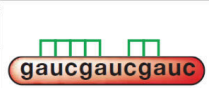
<sup>8</sup>W. Dyrka et al. "Estimating probabilistic context-free grammars for proteins using contact map constraints". *PeerJ* (2019).

# Proteins are 3D objects



Many **crossing** interactions between amino-acids distant in the sequence but close in the structure

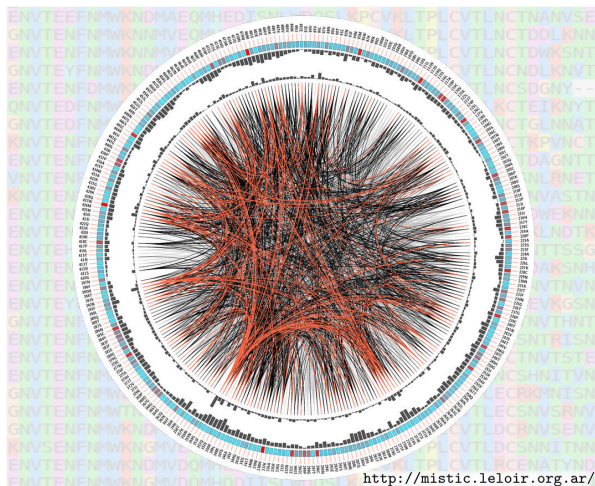
# The Chomsky Hierarchy

Language	Automaton	Grammar	Recognition	Dependencies
Recursively enumerable languages $\cup$	Turing machine 	Unrestricted $Baa \rightarrow A$	Undecidable 	
Context-sensitive languages $\cup$	Linear bounded 	Context sensitive $At \rightarrow aA$	Exponential? 	
Context-free languages $\cup$	Pushdown (stack) 	Context free $S \rightarrow gSc$	Polynomial 	
Regular languages	Finite-state automaton 	Regular $A \rightarrow cA$	Linear 	

Adapted from D. Searls



# Mutual information in HIV-1 gp120 homologs



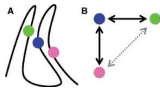
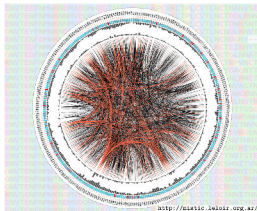
Many (crossing) correlations between MSA columns

Direct Coupling Analysis to the rescue

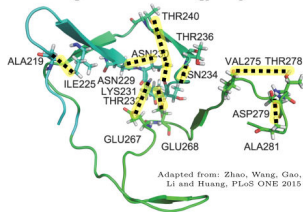


- A recent breakthrough for the prediction of 3D structures by prediction of contacts
- Principle : disentangle direct from indirect effects

Mutual information on HIV-1 gp120 protein



Coevolving residues in HIV-1 gp120 protein



Adapted from: Zhao, Wang, Gao, Li and Huang, PLoS ONE 2015

- Idea: Use DCA for automatic characterization of protein families:
  - Identify important (crossing) dependencies with DCA
  - Build accordingly a syntactic model that can be used in practice. . .

# Choice of DCA method

## CCMpred<sup>9</sup>

- Best one-model precision for contact prediction<sup>10</sup>
- “Structuring” couplings

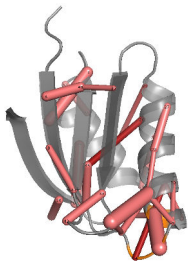


Figure: Top 25 PSICOV predictions

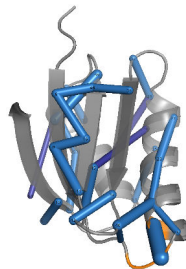


Figure: Top 25 CCMpred predictions

<sup>9</sup>S. Seemayer, M. Gruber, and J. Söding. “CCMpred—fast and precise prediction of protein residue–residue contacts from correlated mutations”. *Bioinformatics* 30.21 (2014), pp. 3128–3130.

<sup>10</sup>S. H. P. de Oliveira, J. Shi, and C. M. Deane. “Comparing co-evolution methods and their application to template-free protein structure prediction”. *Bioinformatics* 33.3 (2017), pp. 373–381.

## 1. Protein sequence query $q$

1CC8:A|PDBID|CHAIN|SEQUENCE MAEIKHYQFNVVMTCSGCSGAVNKVLTKLEPDVSKIDISLEKQLVDVYT

2. Retrieve close homologs and build a MSA (e.g. with HHblits<sup>11</sup>)

```

1CC8: A |PDBID| CHAIN |SEQUENCE
sp|Q54PZ2|ATOX1_DICDI
tr|A7TF58|A7TF58_VANPO
tr|AOA0C7MWI5|AOA0C7MWI5_9SACH
tr|GOWD69|GOWD69_NAUDC
tr|G8ZQK6|G8ZQK6_TORDC
tr|S6E8D5|S6E8D5_ZYGB2
tr|J7R785|J7R785_KAZNA
tr|W1QBQ2|W1QBQ2_DGAPD
tr|H2AUI5|H2AUI5_KAZAF
tr|G8JMM3|G8JMM3_ERECY
tr|S9Q3L9|S9Q3L9_SCHOY
tr|Q01AV4|Q01AV4_OSTTA
tr|E5R4F7|E5R4F7_LEPMJ
tr|R7Z484|R7Z484_CONA1
tr|M3CXY4|M3CXY4_SPHMS
tr|W9XE16|W9XE16_9EURO
tr|Q5BDJ0|Q5BDJ0_EMENI
tr|W3WZP2|W3WZP2_9PEZI
tr|AOA0D2B224|AOA0D2B224_9PEZI
tr|AOA093XHT8|AOA093XHT8_PENMA
tr|AOA074WQB6|AOA074WQB6_9PEZI

```

```

      5    10    15    20    25    30    35    40    45
MAEIKHYQFNVVMTCSGCSGAVNKVLTKLEPDVSKIDISLEKQLVDVYT
...MTYSFFVDMTCGGCSKAVNAILSKIDGVS.NIQIDLENKKVCESS
.STAQHYHFDVVMTCAGCSNAINRVLTRLEPDVSNIEISLEKQTVDDVS
.SNDNHYQFEVVMTCSGCSNAVNKALTRLEPDVSNIDISLENQTVDVHS
.MAENHYQFNVVMTCSGCSNAINRVLTKLEPEVSKIDISLEDQTVDDVT
.SQQNHYQFNVVMSCSGCSNAINKVLSRLEPDVSKIETSLDSQTVDDVY
.MSQNHYHFEVVMSCGCSNAINRVLTKLKPVDVSEIRISLENQTVDDVY
.MSNHYQFDVVMTCASCSNAISKVLTRMEPEVTKFDVSLKQTVDDVQT
.MSAKHYKFDVMTACSGCSNAVNRVLTRL.PGVKNVEISLEKQTVDDVIS
.MIYCYHFNVVMTCSGCSDAIHRSLSKLGPEVTDIDISLENQYVEVFT
.MDTKHYQFQVALACSGCVAAVEKALAKLQPDISKFDISLEKQIVDDVY
...MKYSFNVVMTCDCGCKNAIDRVLNRL.GVDEKEISLEAEVHVTT
.MSTTVTLRCDFACDGCANAVKRILSKDDA...VRTSVEKLVVVV
.MTHTYKFNVTMTCSGCSGAVERVLRKLE.GVESFNVNLETQTAEVVA
.MSEHNYKFNVMSCGGCSGAVERVLRKLD.GVKSFNVSLDQTAEIVA
.MAEHKYKFNVMSCGGCSGAVERVLRKLD.GVKEFNVSLDQTAEITT
.MSEHHYKFNVTMTCSGCSGAVERVLRKLD.GVKNYTVSLDQTAEDVT
.DQEHYKFNVMSCGGCSGAVERVLRKLD.GVKSFVNLDSQTASVVT
.ADNHYYKFNVMSCGGCSGAVDRVLRKLD.GIESYDVSLKQEATVIA
.MSHTYKFNVMSCGGCSGAIIDRVLRKLE.GVDKEYEVSLEKQTAEVHT
.MAEHQYKFNVMSCGGCSGAVERVLRKLDVGVKSYDVSLESQTATVVA
.MSDHTYFNITMTCSGCSGAVERVLRKLD.GVKSFVSLDSQTAFVIT

```

<sup>11</sup>M. Remmert et al. “HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment”. *Nature methods* 9.2 (2012), p. 173.

## 3. Infer a Potts model from MSA



Probability of sequence

$$a = a_1, \dots, a_L$$

$$P(a|w, v) = \frac{1}{Z} \exp \left( \sum_{i=1}^{L-1} \sum_{j=i+1}^L w_{ij}(a_i, a_j) + \sum_{i=1}^L v_i(a_i) \right)$$

Normalization constant

Maximise *pseudo*-likelihood of  $N$  aligned sequences, i.e.:

$$(w, v) = \operatorname{argmax}_{w, v} \sum_{n=1}^N \sum_{i=1}^L \log P(A_i = a_i^n | a_1^n, \dots, a_{i-1}^n, a_{i+1}^n, \dots, a_L^n, v, w)$$

(more tractable and still good precision)

while respecting empirical frequencies:

- $P_i(a) = f_i(a)$
- $P_{ij}(a, b) = f_{ij}(a, b)$



4. Contacts in  $q$  are predicted using Frobenius norm of the couplings

$$\|w_{ij}\| = \sqrt{\sum_a \sum_b w_{ij}(a, b)^2}$$

A larger norm is interpreted as a likelier contact between positions

Using Potts model for homology search

# Using Potts model for homology search

Use whole Potts model  $\mathcal{P}_q$  of  $q$  instead of Frobenius norms

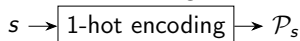
- Use  $\mathcal{P}_q$  to score each possibly homologous sequence  $s$
- Require to compute best alignment of  $s$  in  $\mathcal{P}_q$
- As HHalign for pairs of HMMs<sup>12</sup>, **align directly pairs of Potts models**

↪ A new tool: ComPotts

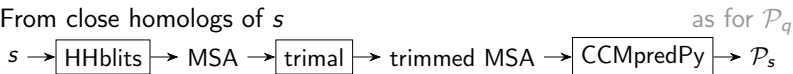


with two options to get Potts model  $\mathcal{P}_s$  of  $s$ :

- One-hot encoding  $v_i(a_i) = 1, w_{ij}(a_i, a_j) = 1$ , others are 0

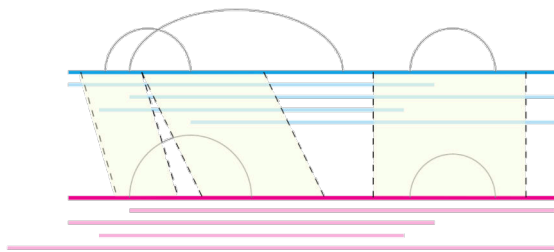


- From close homologs of  $s$



<sup>12</sup>J. Söding. "Protein homology detection by HMM–HMM comparison". *Bioinformatics* 21.7 (2004), pp. 951–960.

# ComPotts (Comparing Potts models)



$$s(A, B) = \sum_{i=1}^{L_A} \sum_{k=1}^{L_B} s_v(v_i^A, v_k^B) x_{ik} + \sum_{i=1}^{L_A-1} \sum_{j=i+1}^{L_A} \sum_{k=1}^{L_B-1} \sum_{l=k+1}^{L_B} s_w(w_{ij}^A, w_{kl}^B) y_{ijkl}$$

$$x_{ik} \geq \sum_{j \in C} y_{ikj} \quad \forall C \in C_{ik}^+, i \in [1, n_A - 1], k \in [1, n_B - 1]$$

$$x_{ik} \geq \sum_{j \in C} y_{jik} \quad \forall C \in C_{ik}^-, i \in [2, n_A], k \in [2, n_B]$$

$$x_{ik} \leq 1 + \sum_{\substack{j \in C \\ s(A_{ij}, B_{kl}) \leq 0}} (y_{ikj} - x_{jl}) \quad \forall C \in C_{ik}^-, i \in [1, n_A - 1], k \in [1, n_B - 1]$$

$$\sum_{i,k \in C} x_{ik} \leq 1 \quad \forall C \in C$$

$$y \geq 0$$

$$x \text{ binary.}$$

- Formulation of Potts model alignment as an Integer Linear Programming (ILP) problem
- Based on Inken Wohlers' solver<sup>13</sup>

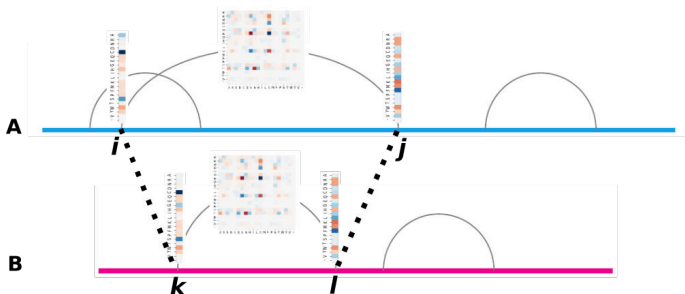
<sup>13</sup>I. Wohlers. "Exact Algorithms For Pairwise Protein Structure Alignment". PhD thesis. Vrije Universiteit, Jan. 2012, pp. 1 –147.

# Scoring alignment of Potts models $A$ and $B$

$$s(A, B) = \sum_{i=1}^{L_A} \sum_{k=1}^{L_B} s_v(v_i^A, v_k^B) x_{ik} + \sum_{i=1}^{L_A-1} \sum_{j=i+1}^{L_A} \sum_{k=1}^{L_B-1} \sum_{l=k+1}^{L_B} s_w(w_{ij}^A, w_{kl}^B) y_{ikjl}$$

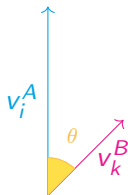
where

- $x_{ik} = 1$  iff position  $i$  of  $A$  and position  $k$  of  $B$  are aligned (otherwise,  $x_{ik} = 0$ )
- $y_{ikjl} = 1$  iff  $x_{ik} = 1$  and  $x_{jl} = 1$  (otherwise,  $y_{ikjl} = 0$ )



# Choice of $s_v(v_i, v_k)$ and $s_w(w_{ij}, w_{kl})$ : scalar products

- $s_v(v_i^A, v_k^B) = \langle v_i^A, v_k^B \rangle$   
→ standard scalar product :  $\langle x, y \rangle = \sum_i x_i y_i$
- $s_w(w_{ij}^A, w_{kl}^B) = \langle w_{ij}^A, w_{kl}^B \rangle_F$   
→ Frobenius scalar product :  $\langle X, Y \rangle_F = \sum_i \sum_j X_{ij} Y_{ij}$



Geometric insight

$$\langle v_i^A, v_k^B \rangle = \|v_i^A\| \|v_k^B\| \cos \theta$$

importance of position  $i$

importance of position  $k$

similarity measure

# Natural extension of the 1D score of a sequence

$$P(a|w, v) = \frac{1}{Z} \exp(\mathcal{H}(a|v, w))$$

$$\mathcal{H}(a|v, w) = \sum_{i=1}^L v_i(a_i) + \sum_{i=1}^{L-1} \sum_{j=i+1}^L w_{ij}(a_i, a_j)$$

$$\sum_{i=1}^L \langle v_i, e_{a_i} \rangle + \sum_{i=1}^{L-1} \sum_{j=i+1}^L \langle w_{ij}, e_{a_i a_j} \rangle F$$

$$e_{a_i a_j} = \begin{pmatrix} 0 & \dots & \dots & 0 & \dots & \dots & 0 \\ \vdots & & & \vdots & & & \vdots \\ & & & 0 & & & \vdots \\ 0 & \dots & 0 & 1 & 0 & \dots & 0 \\ & & & 0 & & & \vdots \\ \vdots & & & \vdots & & & \vdots \\ 0 & \dots & \dots & 0 & \dots & \dots & 0 \end{pmatrix} \leftarrow a_i$$

$a_j$   
↓

$$e_{a_i} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \leftarrow a_i$$

# First experiments

- PDB 1CC8 : Atx1 metallochaperone (*Saccharomyces cerevisiae*)
  - × one homolog  $s$  (150 sequences sampling, identity with 1CC8 : 25%-50%)

One-hot encoding of  $\mathcal{P}_s$

Timeout: 6 hours

	trimmed	not trimmed
$\epsilon = \text{machine epsilon}$	$t \in [11s, 6h]$ , avg: 2h	$t \in [25s, 6h]$ , avg: 4h30
$\epsilon = 1^{14}$	$t \in [8s, 6h]$ , avg: 1h30	$t \in [16s, 6h]$ , avg: 2h

Build  $\mathcal{P}_s$  from homologs of  $s$

Timeout: 6 hours

	trimmed	not trimmed
$\epsilon = \text{machine epsilon}$	$t \in [3s, 6h]$ , avg: 3 min	$t \in [3s, 6h]$ , avg: 8 min
$\epsilon = 1$	$t \in [2s, 6s]$ , avg: 5s	$t \in [3s, 50s]$ , avg: 21s

Tractable time! not for the simpler models?  
Small proteins, easy to align...

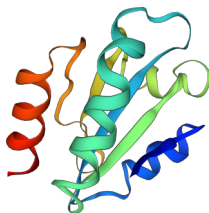
---

<sup>14</sup>  $\simeq$  Energy needed to change one a.a. into another



# Testing the limits on thioredoxins

- Enzymes involved in reduction–oxidation reactions through oxidation of their active site

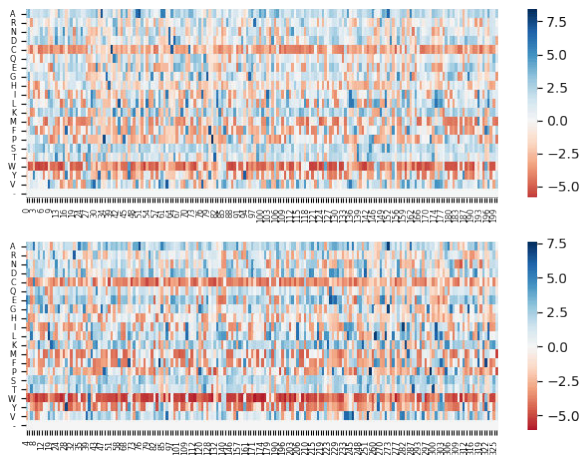


**Figure:** 3D structure of thioredoxin-1 (*Caenorhabditis elegans*) (Q09433)

- 100 amino acids on average
- Between 15 and 20% sequence identity within the family
  - known to be hard to align

# An example of failure

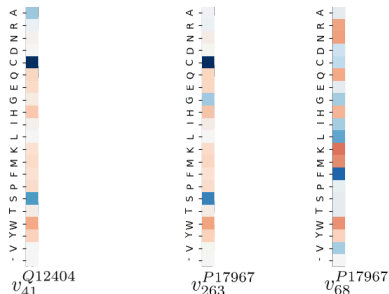
$\{v_i^{Q12404}\}_i$  and  $\{v_i^{P17967}\}_i$  aligned by ComPotts:



Even well-conserved positions of the active site are not aligned

# The trouble with scalar product alone

- A well-conserved column  $i$  may have a smaller  $\|v_i\|$  than a less conserved column  $j$



$$\|v_{262}^{P17967}\| \simeq 10.5 < \|v_{68}^{P17967}\| \simeq 11.4$$

- It may be more profitable to align many less conserved columns than to align fewer well-conserved columns with each other

# An idea

- Use rescaling function :  $f(x) = \text{sign}(x)(e^{|x|} - 1)$

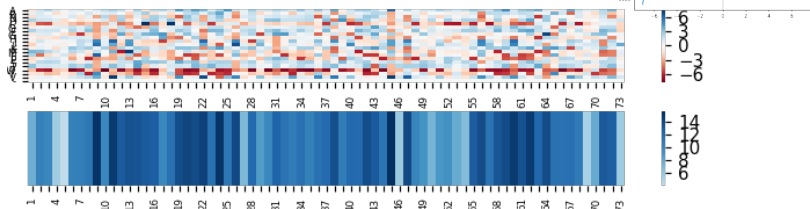


Figure: Before rescaling

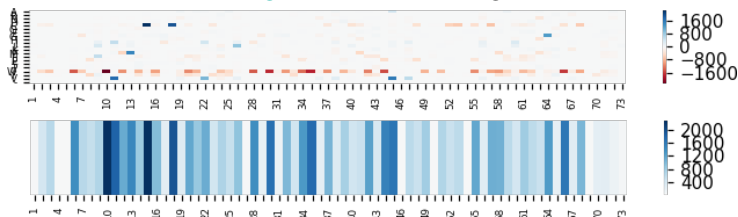
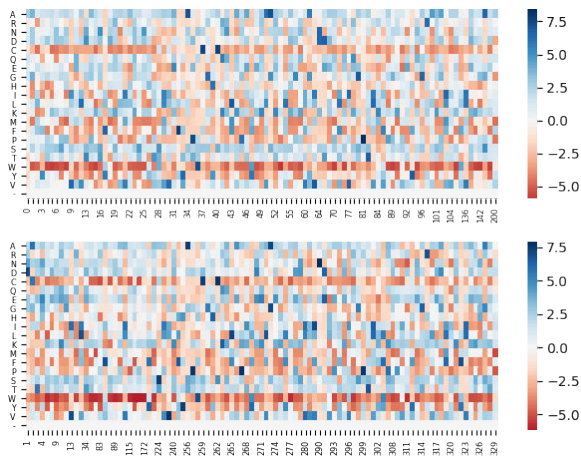


Figure: After rescaling

It's better :-)

$\{v_i^{Q12404}\}_i$  and  $\{v_i^{P17967}\}_i$  aligned by ComPotts (wo couplings!)



# To be continued...

- How to rescale also consistently the couplings  $w_{ij}$ ?
  - Slightly change the rescaling function  $f(x) = \text{sign}(x)(\beta e^{\alpha|x|} - \gamma)$ ?
- Other similarity functions?...
- Introduce gap costs
- Constrain Potts model inference?
  - Canonical Potts model?
  - Better control amplitude of vectors and matrices?

## Conclusion so far...

- Good news: alignment to Potts model is tractable
- A surprise: may require a transformation to a Potts model
- A working efficient implementation
- Quality of alignment can still be improved...

Thanks for your attention!  
Ideas, remarks, suggestions are welcome.  
See you next to our poster. . .



# Bibliography I

- [Alt+90] S. F. Altschul et al. “Basic local alignment search tool”. *Journal of molecular biology* (1990).
- [Edd98] S. R. Eddy. “Profile hidden Markov models.”. *Bioinformatics* 14.9 (1998), pp. 755–763.
- [Ste+19] M. Steinegger et al. “HH-suite3 for fast remote homology detection and deep protein annotation”. *bioRxiv* (2019), p. 560029.
- [Ker08] G. Kerbellec. “Apprentissage d’automates modélisant des familles de séquences protéiques”. PhD thesis. Université de Rennes 1, Apr. 2008, p. 139.
- [Bre+13] A. Bretaudeau et al. “CyanoLyase: a database of phycobilin lyase sequences, motifs and functions”. *Nucleic Acids Research* 41.Database-Issue (2013), pp. 396–401.

- [CGN14] F. Coste, G. Garet, and J. Nicolas. “A bottom-up efficient algorithm learning substitutable languages from positive examples”. *ICGI*. 2014.
- [Dyr+19] W. Dyrka et al. “Estimating probabilistic context-free grammars for proteins using contact map constraints”. *PeerJ* (2019).
- [SGS14] S. Seemayer, M. Gruber, and J. Söding. “CCMpred—fast and precise prediction of protein residue–residue contacts from correlated mutations”. *Bioinformatics* 30.21 (2014), pp. 3128–3130.
- [OSD17] S. H. P. de Oliveira, J. Shi, and C. M. Deane. “Comparing co-evolution methods and their application to template-free protein structure prediction”. *Bioinformatics* 33.3 (2017), pp. 373–381.

- [Jon+11] D. T. Jones et al. “PSICOV: precise structural contact prediction using sparse inverse covariance estimation on large multiple sequence alignments”. *Bioinformatics* 28.2 (2011), pp. 184–190.
- [Rem+12] M. Remmert et al. “HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment”. *Nature methods* 9.2 (2012), p. 173.
- [Söd04] J. Söding. “Protein homology detection by HMM–HMM comparison”. *Bioinformatics* 21.7 (2004), pp. 951–960.
- [Woh12] I. Wohlers. “Exact Algorithms For Pairwise Protein Structure Alignment”. PhD thesis. Vrije Universiteit, Jan. 2012, pp. 1–147.