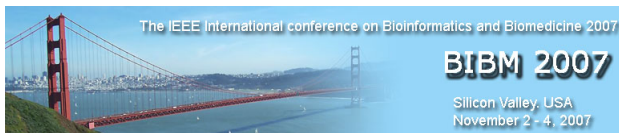# A Protocol to Detect Local Affinities Involved in Proteins Distant Interactions

Christophe N. Magnan, Cécile Capponi, François Denis
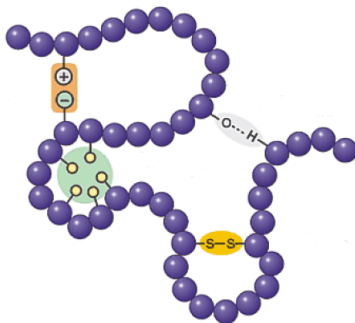
LIF, CNRS, France

The IEEE International conference on Bioinformatics and Biomedicine 2007

BIBM 2007

Silicon Valley, USA
November 2 - 4, 2007

**Introduction**
Revealing local affinities
Experimentations
Conclusions

**Proteins distant interactions**
Disulfide bridges example
Local information

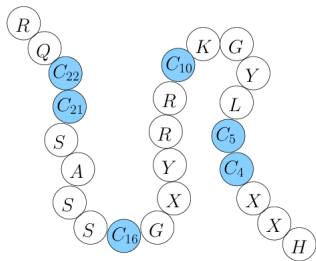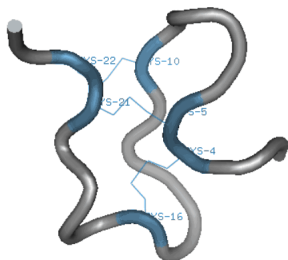# Proteins distant interactions - examples

Ionic links, hydrogen bonds, hydrophilic interactions, salt bridges, disulfide bridges, Van Der Waals forces, ...



3D structure constrained/stabilized by these interactions

**Introduction**
Revealing local affinities
Experimentations
Conclusions

Proteins distant interactions
**Disulfide bridges example**
Local information

# Prediction of disulfide bridges: a two-stage process

Introduction
Revealing local affinities
Experimentations
Conclusions

Proteins distant interactions
Disulfide bridges example
Local information

# Prediction of disulfide bridges: a two-stage process

**Introduction**
Revealing local affinities
Experimentations
Conclusions

Proteins distant interactions
**Disulfide bridges example**
Local information

# Prediction of disulfide bridges: a two-stage process

**Introduction**
Revealing local affinities
Experimentations
Conclusions

Proteins distant interactions
Disulfide bridges example
**Local information**

# Local environments of bonded amino-acids

Introduction
Revealing local affinities
Experimentations
Conclusions

Proteins distant interactions
Disulfide bridges example
Local information

# Local environments of bonded amino-acids



Is there information carried by local environments involved in the formation of bonds such as disulfide bridges?

**Introduction**
Revealing local affinities
Experimentations
Conclusions

Proteins distant interactions
Disulfide bridges example
**Local information**

## Local information?

Are the local environments involved in interactions?

- $\beta$-sheets: there is local information
- Disulfide/salt bridges: no biological evidence
- Some biologists and biochemists skeptical

**Introduction**
Revealing local affinities
Experimentations
Conclusions

Proteins distant interactions
Disulfide bridges example
**Local information**

# Local information?

Are the local environments involved in interactions?

- $\beta$-sheets: there is local information
- Disulfide/salt bridges: no biological evidence
- Some biologists and biochemists skeptical
- Always used to predict disulfide bridges

**Introduction**
Revealing local affinities
Experimentations
Conclusions

Proteins distant interactions
Disulfide bridges example
**Local information**

## Local information?

Are the local environments involved in interactions?

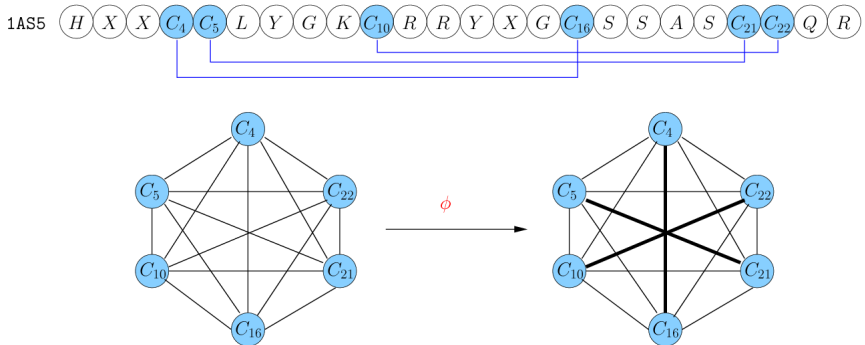- $\beta$-sheets: there is local information
- Disulfide/salt bridges: no biological evidence
- Some biologists and biochemists skeptical
- Always used to predict disulfide bridges

- Is it possible to detect such information?
- Is it possible to show that there exists an **affinity** between local environments of bonded residues involved in the pairing of these residues?

Introduction
**Revealing local affinities**
Experimentations
Conclusions

**Modelling the data**
A first approach
A reasonable approach

## Model

- $\Sigma$: the set of 20 amino-acids,
- $\mathcal{P} \subset \Sigma^*$: proteins containing an even number of amino-acids involved in bridges
- $\mathcal{P}_l \subset \mathcal{P}$, proteins with $2l$ amino-acids involved in bridges
- $\phi$, a function which associates the correct connectivity to a protein in $\mathcal{P}$

Introduction
**Revealing local affinities**
Experimentations
Conclusions

**Modelling the data**
A first approach
A reasonable approach

# Model



The prediction of interactions between amino-acids amounts to approximating $\phi$ with the highest precision

Introduction
**Revealing local affinities**
Experimentations
Conclusions

**Modelling the data**
A first approach
A reasonable approach

## Model

Local environments: segments centered on bonded amino-acids of size $2r + 1$ are considered.

- $P$ a distribution over $\mathcal{P}$
- $\Omega_r = \Sigma^{2r+1}$ the set of proteins segments of size $2r + 1$

- For $w, w' \in \Omega_r$ let:
  - $P(w)$ the probability that $w$ is a local environment
  - $P(B(w, w')|w, w', l)$ the probability that $w$ and $w'$ are bonded knowing that there are distinct local environments of amino-acids involved in interactions into a protein $p \in \mathcal{P}_l$

Introduction
**Revealing local affinities**
Experimentations
Conclusions

Modelling the data
**A first approach**
A reasonable approach

## Local information?

Let $p$ be a protein with $l$ bridges ($2l$ involved amino-acids).

$P(B(w, w')|w, w', l) = \dfrac{1}{2l - 1} \Leftrightarrow$ No local information for pairing amino acids

- a probabilistic way to determine if the local context of bonded residues is involved into the formation of the bridges

- but, estimating directly these probabilities is impossible:

  $r = 3 \rightarrow |\{(w, w'), w, w' \in \Omega_r\}| = 20^{12} \simeq 4 \cdot 10^{15}$, while only few hundreds examples are available in databases!

Introduction
**Revealing local affinities**
Experimentations
Conclusions

Modelling the data
A first approach
**A reasonable approach**

# An affinity function $g$

The solution we propose:

To suppose the existence of an affinity function $g : \Omega_r^2 \rightarrow Y$ ($|Y|$ small) such that:

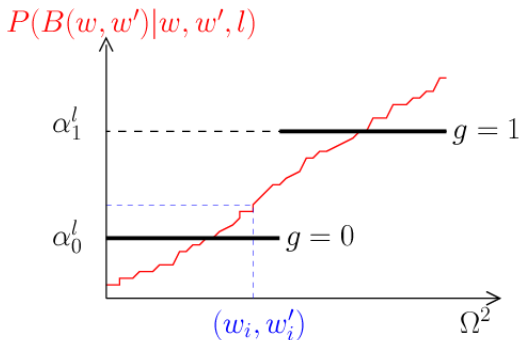$$g(w_1, w_2) = g(w_1', w_2') \Rightarrow P(B(w_1, w_2)|w_1, w_2, I) \simeq P(B(w_1', w_2')|w_1', w_2', I)$$

and

$$y < y' \Rightarrow P(B(w_1, w_2)|g(w_1, w_2) = y) < P(B(w_1', w_2')|g(w_1', w_2') = y')$$

Introduction
**Revealing local affinities**
Experimentations
Conclusions

Modelling the data
A first approach
**A reasonable approach**

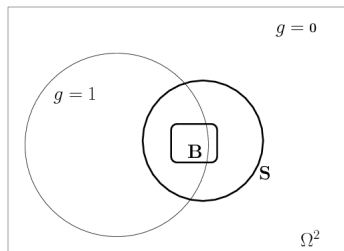# A simple case: $Y = \{0, 1\}$

With $Y = \{0, 1\}$, pairs of local environments are partitioned into two classes, corresponding to two affinity levels and:

$$P(B(w, w')|w, w', l) \simeq P(B(w, w')|g(w, w'), l) = \left\{ \begin{array}{l} \alpha_1^l \ \text{if} \ g(w, w') = 1 \\ \alpha_0^l \ \text{if} \ g(w, w') = 0 \end{array} \right.$$

Introduction
**Revealing local affinities**
Experimentations
Conclusions

Modelling the data
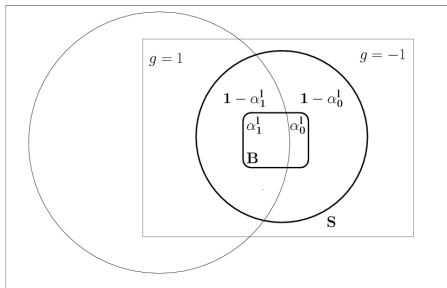A first approach
**A reasonable approach**

## Observations as indirect information on g

The observed classes (bonded or non-bonded) of examples issued from experiments do not carry direct information about g.

Introduction
**Revealing local affinities**
Experimentations
Conclusions

Modelling the data
A first approach
**A reasonable approach**

## Observed pairs as noisy examples of g

The pairs such that

- $g=1$ correspond to observing a bridge with noise $\eta^+=1-\alpha_1^l$
- $g=0$ correspond to non-bonded pairs with noise $\eta^-=\alpha_0^l$



- generalization of the *uniform classification noise* ($\eta^+ = \eta^-$)
- referred to as *class-conditional classification noise* (CCCN)

Introduction
**Revealing local affinities**
Experimentations
Conclusions

Modelling the data
A first approach
**A reasonable approach**

# Setting up the protocol to learn $g$

**If** a local information exists

**If** it can be represented by a function learnable under CCCN
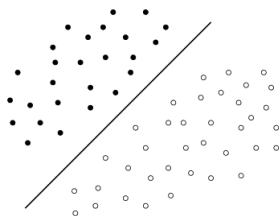
**then** we should be able to detect, extract and evaluate it

assuming that we have access to a sufficient number of examples

Introduction
Revealing local affinities
**Experimentations**
Conclusions

**CCCN-algorithms**
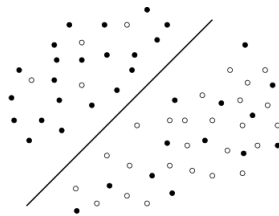Experimentation protocol
Results

What can we learn under CCCN?

- some theoretical results
- they can not be used in practice
- methods such as Soft-margins SVM cannot handle data corrupted by CCCN
- New methods have to be created

Introduction
Revealing local affinities
**Experimentations**
Conclusions

**CCCN-algorithms**
Experimentation protocol
Results

## Perceptron CCCN

- we propose an algorithm to learn *linear threshold functions* from examples corrupted by CCCN
- a generalization of the Perceptron algorithm



No noise                    CCCN

Introduction
Revealing local affinities
**Experimentations**
Conclusions

CCCN-algorithms
**Experimentation protocol**
Results

## Datasets

- 1 dataset of proteins featuring salt bridges: G3D
    - 1836 internal salt bridges in 570 proteins
    - created from PDB by Christophe Geourjon (IBCP, Lyon, France) in 2005

- 1 dataset of proteins featuring disulfide bridges: SPX
    - 1676 internal disulfide bridges within 567 proteins
    - created from Swiss-Prot by Jianlin Cheng and Pierre Baldi (Irvine, California) in 2005

- proteins containing from 2 to 5 bonds

Introduction
Revealing local affinities
**Experimentations**
Conclusions

CCCN-algorithms
**Experimentation protocol**
Results

## Coding of local environments pairs

For a protein containing $l$ bridges:

- $l(2l - 1)$ pairs of local environments
- radius $r = 6$ ($|w| = 13$)
- each local environments pair $(w, w')$ is described as follows:
  - 169 amino-acids pairs $(a_i, a_j)$, with $a_i \in w$ and $a_j \in w'$ ($i, j \in \{1, ..., 13\}$)

  - $(w, w')$ is modeled with a vector of $\mathbb{R}^m$ with:
    - $m$ is the number of ordered pairs of amino-acids in $\Sigma$ ($m = 231$)
    - each coordinate is the number of time the corresponding pair is observed in $(w, w')$

Introduction
Revealing local affinities
**Experimentations**
Conclusions

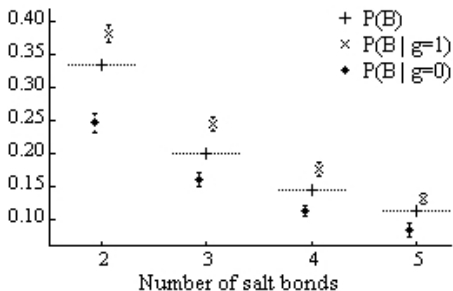CCCN-algorithms
**Experimentation protocol**
Results

## Experiments and studied criteria

- we launch 5 10-fold cross-validations for both kinds of bonds
- two criteria are studied:
  - $P(B|g = 1)$, the probability to observe a bond knowing that the pair is predicted to have a high level of affinity
  - $P(B|g = 0)$, the probability to observe a bond knowing that the pair is predicted to have a low level of affinity

Introduction
Revealing local affinities
**Experimentations**
Conclusions

CCCN-algorithms
Experimentation protocol
**Results**

## Salt bridges

A clear signal is detected:
$\forall l \in \{2, 3, 4, 5\},\ P(B|g = 1, l) > P(B|g = 0, l)$

Introduction
Revealing local affinities
**Experimentations**
Conclusions

CCCN-algorithms
Experimentation protocol
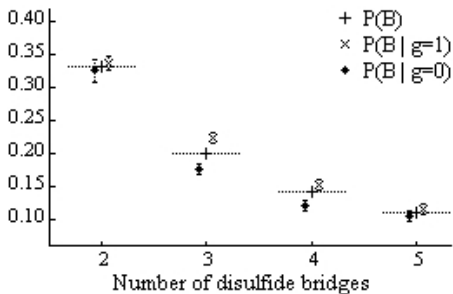**Results**

## Salt bridges

The detected affinities might be explained either by

- the ionic nature of salt bridges
- the hydrophilic property of many residues around salt bridges

Introduction
Revealing local affinities
**Experimentations**
Conclusions

CCCN-algorithms
Experimentation protocol
**Results**

## Disulfide bridges

Results are not as clear as expected:
$\forall l \in \{2, 3, 4, 5\},\ P(B|g=1, l) \simeq P(B|g=0, l)$

Introduction
Revealing local affinities
**Experimentations**
Conclusions

CCCN-algorithms
Experimentation protocol
**Results**

## Disulfide bridges

These results may be explained by several independent reasons:

- Biology reality: there might be no local information that would guide the formation of disulfide bridges
- Learning a function in an unsuitable function class: the function $g$ that we try to learn might be not representable by a linear threshold function.
- ...

This work give us no hint on which assumption is the most probable

## Conclusions

- a machine-learning based protocol to answer the question of the presence of local affinities
- independent from the contact
- results on salt bridges validate this protocol
- disulfide bridges remain an open question