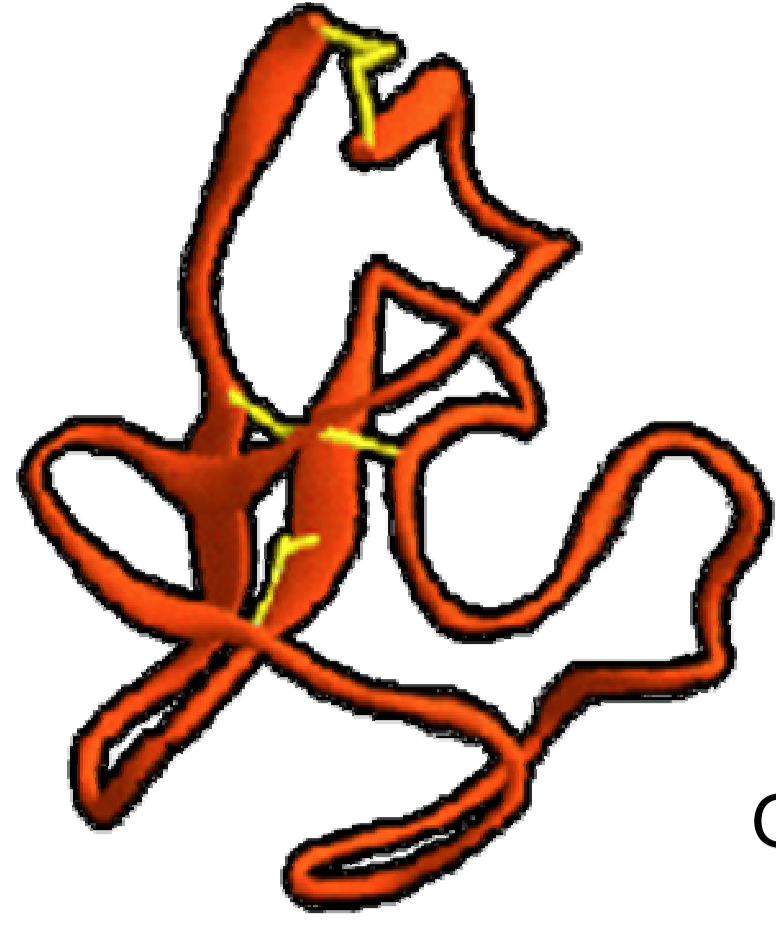


Projet GenoTo3D

Apprentissage automatique appliqué à la prédiction de la structure tertiaire des protéines



Guermeur Y¹, Benabdelsem K², Bréhélin L³, Capponi C⁴, Coste F⁶, Darcy Y¹, Deléage G², Denis F⁴, Gascuel O³, Geourjon C², Gibrat JF⁵, Jacquemin I⁶, Magnan C⁴, Marin A⁵, Martin J⁵, Monfrini E¹, Nicolas J⁶, Ralaivola L⁴, Taly JF⁵
1 : LORIA-Nancy, 2 : IBCP-Lyon, 3 : LIRMM-Montpellier, 4 : LIF-Marseille, 5 : MIG-Jouy en Josas, 6 : IRISA-Rennes

Prédiction de la structure et apprentissage automatique

Problème d'apprentissage sur de grandes quantités de données
Contexte biologique Exploitation fonctionnelle des informations provenant des grands programmes de séquençage des génomes : passe par la connaissance de la structure 3D des protéines. C'est cette structure 3D qui est responsable de la fonction biologique.

1. Arrivée massive de séquences protéiques (croissance exponentielle des bases)

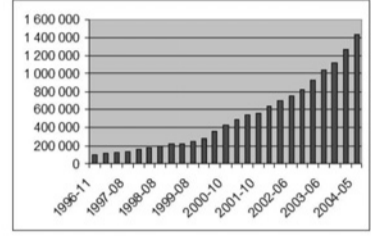


FIG. 1 - Croissance de la banque internationale TREMBL de 1996 à 2004

2. Détermination expérimentale de la structure 3D : tâche très lourde... lorsqu'elle est réalisable
=> Nécessité de passer d'une approche biochimique à une approche prédictive

Problème central en biologie permettant d'aborder l'essentiel des grandes questions ouvertes en traitement de données séquentielles

Différents niveaux d'organisation structurale des protéines

• Séquence ou structure primaire (1 536 117 séquences communes)
MEELKAKKIFVVGQSGKGTQCEKIVKRYGTHLSTC...

• Structure secondaire

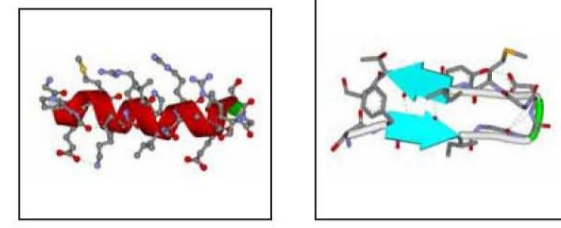


FIG. 2 - Elements structuraux périodiques : hélice α (à gauche) et brins β (à droite)

• Structure tertiaire (27 112 structures 3D communes)



Projet GENOTO3D de l'ACI "Masses de Données"

• Apprentissage automatique appliqué à la prédiction de la structure tertiaire des protéines

"L'objet de ce projet est, dans le contexte de la prédiction de la structure tertiaire des protéines, de mettre en évidence des problèmes de prédiction sur des séquences génériques et difficiles, et de proposer des méthodes susceptibles de faire progresser l'état de l'art dans le domaine."

• Participants

1. Projet MODBIO du LORIA (Nancy)
2. Laboratoire de Bioinformatique et RMN Structurales (LBRS) de l'IBCP (Lyon)
3. Equipe "Bases de Données et Apprentissage Automatique" du LIF (Marseille)
4. Projet Symbiose de l'IRISA (Rennes)
5. Equipe "Méthodes et algorithmes pour la bioinformatique" (MAB) du LIRMM (Montpellier)
6. Unité Mathématique, Informatique et Génome (MIG) de l'INRA, centre de Jouy-en-Josas

Approche modulaire et hiérarchique

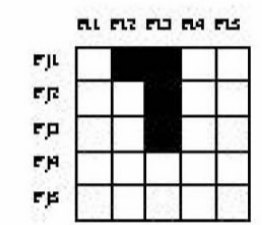
• Un ensemble de sous-problèmes et de reformulations du problème

- Prédiction des ponts disulfures et des ponts salins : IBCP, IRISA, LIF, LORIA
- Prédiction de la structure secondaire (feuilles β , ...) : LIRMM, LORIA, MIG
- Prédiction par homologie ou analogie et reconnaissance des coeurs structuraux : IBCP
- Prédiction par *threading* : IRISA, MIG
- Prédiction *ab initio* (de novo) : MIG

Prédiction des ponts disulfures

Modélisation probabiliste des appariements d'acides aminés autour des cystéines

• Extraction de la PDB des cartes de contact de chaque pont disulfure



L'extraction est paramétrée par la taille K des fenêtres et la distance minimale d entre AA sous laquelle on considère qu'il y a un contact

• Modèle probabiliste des contacts :

$\theta(k, l)$ est la probabilité d'un contact à la position (k, l)

• blanc : la probabilité du contact à cette position est élevée ;
• rouge : la probabilité est faible.

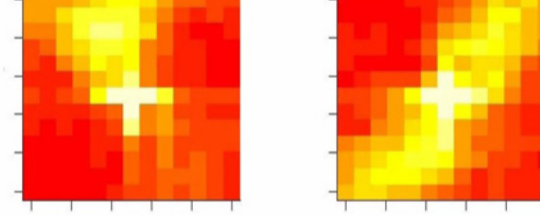
On peut calculer la probabilité d'une carte sous ce modèle :

$$P(F, F') = \prod_{i,j} \theta(i, j)^{I_{i,j}} \prod_{i,j} (1 - \theta(i, j))^{1 - I_{i,j}}$$

• Les cartes sont utilisées pour apprendre :

- une distribution de paires d'AA pour les contacts et les non contacts : $P(a, b | \text{contact})$ et $P(a, b | \text{non contact})$;
- un modèle de mélange de C classes (grâce à un algorithme de type EM).

Exemple de modèle de mélange à 2 classes :



• La probabilité d'observer deux fenêtres d'AA en contact est alors :

$$P(F, F') = \sum_{i,j} P(i, j) P(F, F' | i, j)$$

• Etant donnée une protéine à 2N cystéines, on peut utiliser cette modélisation pour prédire les N ponts disulfures les plus probables.

Prédiction des ponts disulfures par PLI

Knowledge
properties(A, hydrophobic, small, tiny) ...
type(-7, far), type(-7, left) ...
patterns(quarter(A,B,C,D), [L,A,L,B,L,C,L,D]) ...
properties(A, LA, property(B, LB), property(C, LC), property(D, LD)) ...

Positive
exemple(context([L, G, V, A, A, A, L, L, L, P, C, H]))
exemple(context([L, C, A, L, L, H, P, L, G, A, L, M, A, Q]))

Negative
exemple(context([L, A, K, V, G, G, A, C, T, P, V, A, F, D]))
exemple(context([L, H, M, E, E, D, P, C, K, L, V, K, K]))

Background Knowledge
-modele(+context(+context))?
-modele(+pattern(+context,+properties)?) ...
-set(ouire, SP) ...
prune(+pattern(+context,+properties)) ...

properties(A, hydrophobic, small, tiny) ...
type(-7, far), type(-7, left) ...
patterns(quarter(A,B,C,D), [L,A,L,B,L,C,L,D]) ...
properties(A, LA, property(B, LB), property(C, LC), property(D, LD)) ...
exemple(context([L, G, V, A, A, A, L, L, L, P, C, H])) ...
exemple(context([L, C, A, L, L, H, P, L, G, A, L, M, A, Q])) ...

PROGOL

Results
Exemple(A) :-
pattern(A, quarter(A,B,C,D), [L,A,L,B,L,C,L,D]),
properties(A, LA, property(B, LB), property(C, LC), property(D, LD)).

Exemple(A) :-
pattern(A, quarter(A,B,C,D), [L,A,L,B,L,C,L,D]),
properties(A, LA, property(B, LB), property(C, LC), property(D, LD)).



Protocole pour détecter la présence d'information locale pour la prédiction de contacts entre acides aminés

Motivations

Les ponts disulfures sont des liaisons covalentes entre cystéines oxydées qui forment un élément du repliement d'une protéine. Mais il n'est pas clair si les ponts participent au repliement ou en sont la cause. En particulier : les voisins des cystéines oxydées sur les séquences primaires contiennent-ils une information sur les ponts qu'elles forment ?

Problématique : poser un cadre formel pour détecter la présence d'une information au voisinage de deux acides aminés permettant de (contribuer à) prédire s'ils sont en contact ou non.

Modélisation (définie dans le cadre de la prédiction des ponts disulfures)

Protéines : (Σ, P) où $\Sigma = \{\text{acides aminés}\}$ et P dist. de probabilité. $P \in \Sigma^2$: protéines contenant un nombre pair de cystéines oxydées.

Pour $w, w' \in \Omega = \{w \in \Sigma^{2n} | |w|_r = 1\}$ (cystéines) (contextes locaux centrés sur une cystéine), on définit

$P(B(w, w') | C(w, w', l))$: probabilité que w et w' forment un pont sachant que ce sont des contextes locaux distincts de cystéines oxydées d'une protéine en contenant $2l$.

Absence d'information locale
 $\forall w, w' \in \Omega, \forall l, P(B(w, w') | C(w, w', l)) = 1/(2l - 1)$.

Inévitablement d'estimer directement $P(B(w, w') | C(w, w', l))$. Idée : supposer l'existence d'une fonction d'affinité $g : \Omega \rightarrow Y$ (Y petit tq $g(w_1, w_2) = g(w'_1, w'_2) \Rightarrow \forall l, P(B(w_1, w_2) | C(w_1, w_2, l)) = P(B(w'_1, w'_2) | C(w'_1, w'_2, l))$.

Cas le plus simple : $Y = \{0, 1\}$. Les paires de fenêtres se répartissent en deux classes, correspondant à deux niveaux d'affinité : faible (0) et fort (1). On devra alors pouvoir observer

$P(B(w, w') | g(w, w')) = 1, l) > P(B(w, w') | g(w, w') = 0, l)$.

Affinité, ponts et bruit de classification
Sous l'hypothèse de l'existence d'une telle fonction d'affinité g ,

$P(B(w, w') | C(w, w', l)) = P(B(w, w') | g(w, w'), l) = \begin{cases} 1 & \text{si } g(w, w') = 1 \\ 0 & \text{si } g(w, w') = 0 \end{cases}$

L'observation d'un pont (H) correspond à $g = 1$ avec un bruit de classification différent de paramètres : $\eta^1 = 1 - \alpha^1$ et $\eta^0 = \alpha^0$.

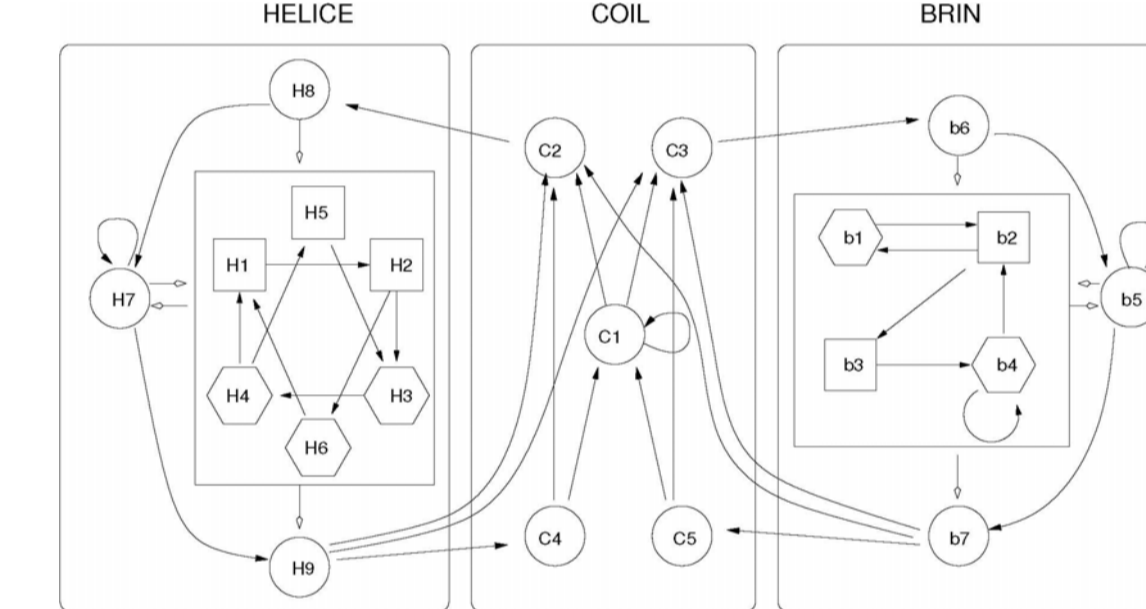
Prédiction de la structure locale des protéines

Notre objectif est de prédire la structure locale en terme de structure secondaires (hélices α , brins β , boucles) et de zones d'angles dièdres Phi/Psi qui apportent une information sur la structure des boucles.

La structure locale de la protéine constitue le processus caché du modèle de chaîne de Markov caché (HMM) avec une mémoire d'ordre 1. La séquence de la protéine est le processus observé. Les acides aminés sont émis indépendamment conditionnellement à la structure locale.

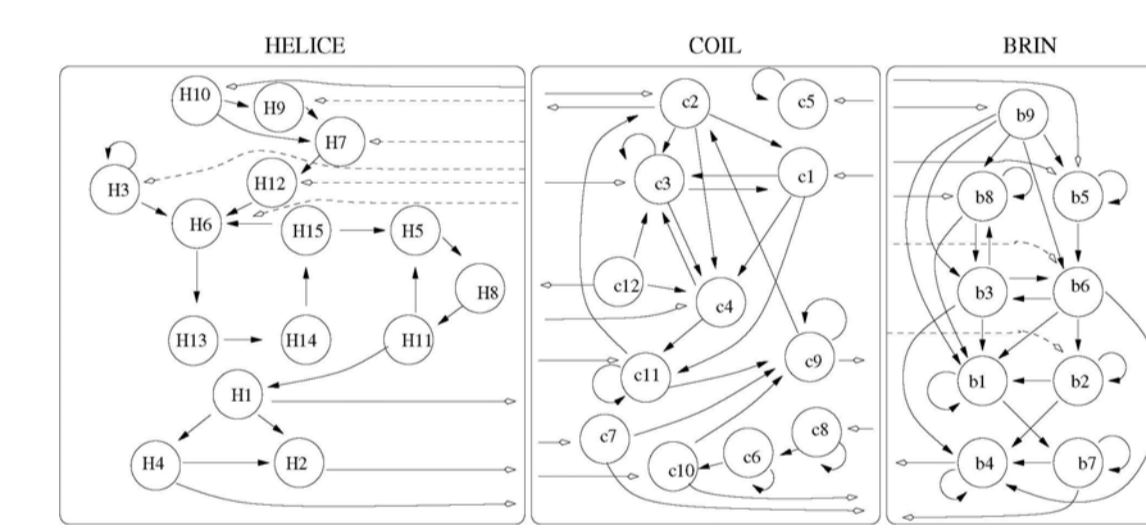
Chaque classe structurale est modélisée par un ensemble d'états cachés. Deux stratégies ont été mises en place pour construire ces modèles :

➤ Modélisation des structures secondaires à partir d'a priori biologiques (modélisation des hélices amphiphiles) et d'une étude des mots exceptionnels dans les brins.



HMM à 21 états cachés. Les états carrés privilégient les résidus polaires, les états hexagonaux les résidus hydrophobes. Le taux de bonne prédiction obtenu est de 65%.

➤ Choix du nombre d'états cachés selon des critères de performance et le critère BIC

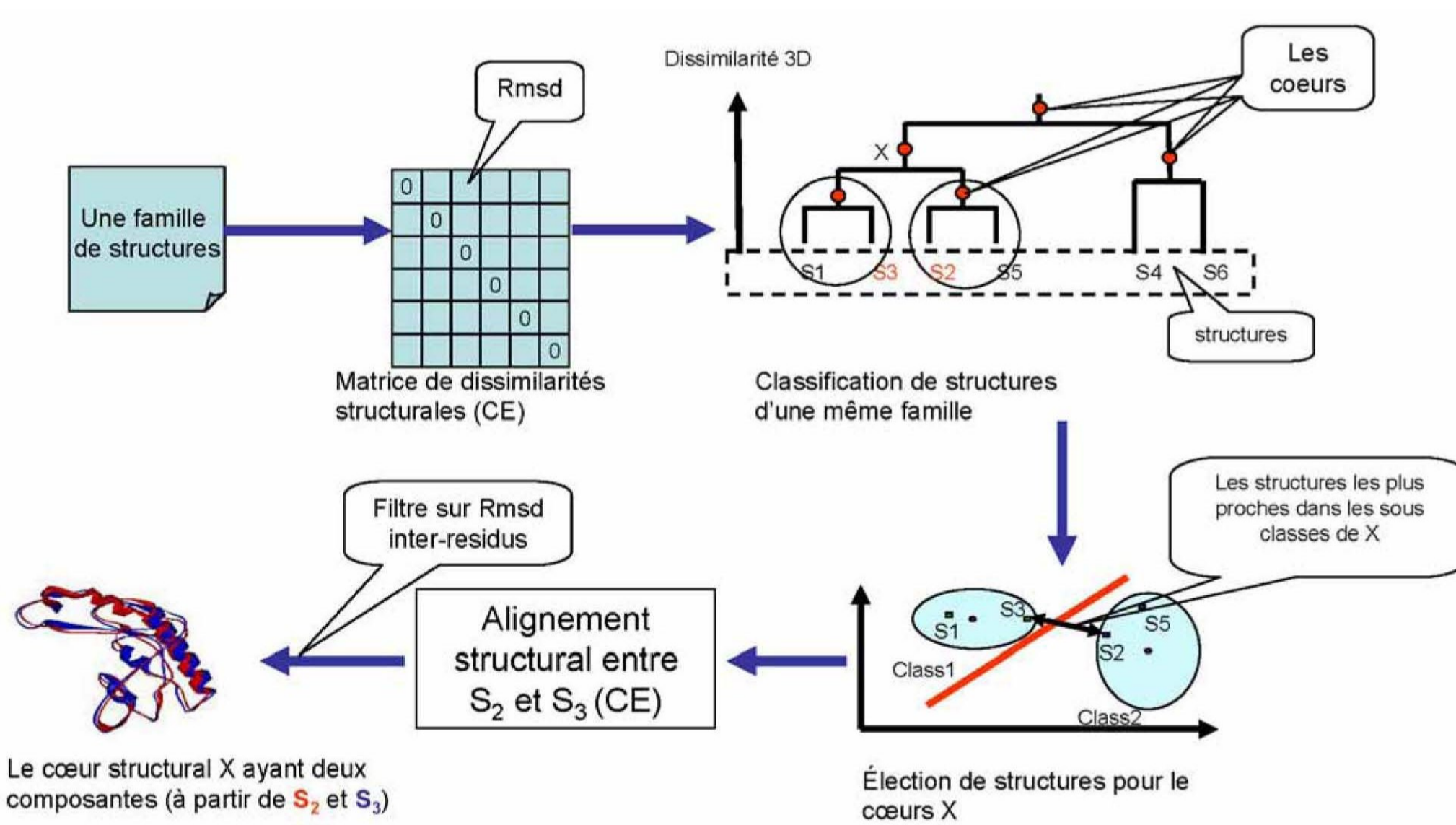


HMM à 36 états cachés. Seules les transitions les plus probables sont indiquées. Le taux de bonne prédiction obtenue est de 68%.

La prédiction est améliorée grâce aux séquences homologues. Les contributions des séquences homologues sont combinées avec les pondérations de Henikoff. Le taux de bonne prédiction atteint ainsi 76%. La même méthodologie appliquées aux zones d'angles permet d'atteindre un taux de prédiction de 78%.

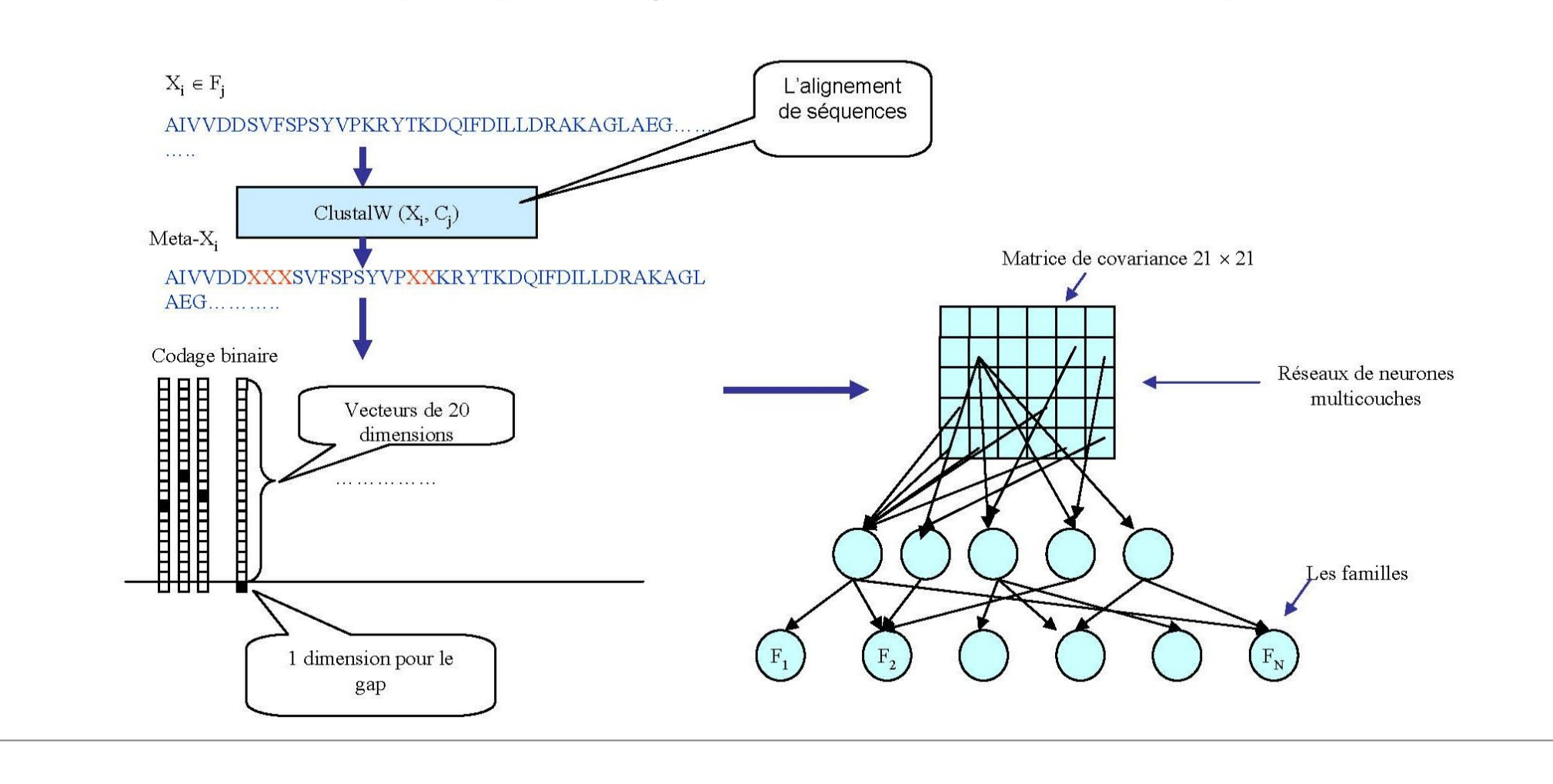
Apprentissage à partir des coeurs structuraux

Alignement et classification de structures pour l'extraction des coeurs structuraux

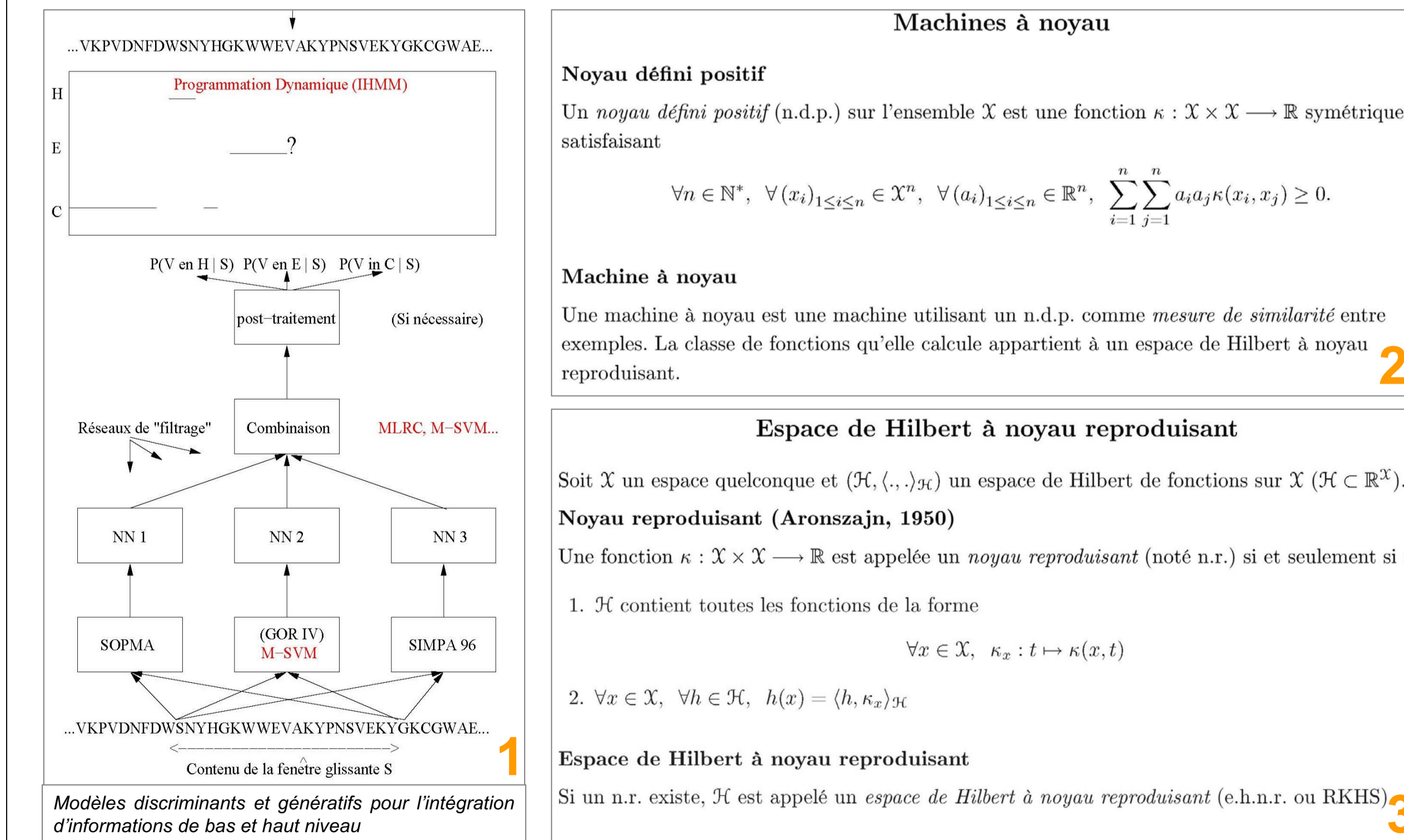


Codage matriciel et modélisation connexionniste

1. Apprentissage : alignement de chaque séquence avec le cœur structural extrait de sa famille
2. Test : alignement de chaque séquence avec chacun des coeurs structuraux déterminés
3. Modélisation de chaque séquence alignée par une matrice de covariance présentée à un PMC



Approche hiérarchique de la prédiction de la structure secondaire



Machines à noyau

Noyau défini positif
Un noyau défini positif (n.d.p.) sur l'ensemble X est une fonction $\kappa : X \times X \rightarrow \mathbb{R}$ symétrique satisfaisant

$$\forall n \in \mathbb{N}^+, \forall (x_i)_{1 \leq i \leq n} \in X^n, \forall (a_i)_{1 \leq i \leq n} \in \mathbb{R}^n, \sum_{i,j=1}^n a_i a_j \kappa(x_i, x_j) \geq 0.$$

Machine à noyau
Une machine à noyau est une machine utilisant un n.d.p. comme mesure de similarité entre exemples. La classe de fonctions qu'elle calcule appartient à un espace de Hilbert à noyau reproduisant.

Espace de Hilbert à noyau reproduisant

Soit X un espace quelconque et $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ un espace de Hilbert de fonctions sur X ($\mathcal{H} \subset \mathbb{R}^X$).
Noyau reproduisant (Aronszajn, 1950)
Une fonction $\kappa : X \times X \rightarrow \mathbb{R}$ est appelée un noyau reproduisant (noté n.r.) si et seulement si :

1. \mathcal{H} contient toutes les fonctions de la forme $\forall x \in X, \kappa_x : t \mapsto \kappa(x, t)$
2. $\forall x \in X, \forall h \in \mathcal{H}, h(x) = \langle h, \kappa_x \rangle_{\mathcal{H}}$

Espace de Hilbert à noyau reproduisant
Si un n.r. existe, \mathcal{H} est appelé un espace de Hilbert à noyau reproduisant (e.h.n.r. ou RKHS).

Nouveau noyau pour le traitement des séquences protéiques

Expression analytique (structure primaire seule)
 (x, x') : vecteurs des descriptions de deux contenus de fenêtres (segments) à comparer

$$\kappa_{\mu, \theta, D}(x, x') = \exp \left(-\mu \sum_{i=-n}^n \theta_i^2 \|x_i - x'_i\|^2 \right)$$

avec $D = (d_{jk})_{1 \leq j, k \leq 22}, \|x_i - x'_i\|^2 = d_{jj} + d_{kk} - 2d_{jk}$ pour $x_i = a_j$ et $x'_i = a_k$

Extension pour le traitement des alignements multiples
Immédiate : x est remplacé par \tilde{x} tel que $\tilde{x}_i = \sum_{j=1}^{22} \theta_{ij} a_j$

$$\sum_{j=1}^{22} \theta_{ij} a_j, \sum_{k=1}^{22} \theta'_{ik} a_k = \sum_{j=1}^{22} \sum_{k=1}^{22} \theta_{ij} \theta'_{ik} (a_j, a_k)$$

Références

- Ingrid Jacquemin. Découverte de motifs relationnels en bioinformatique : application à la prédiction des ponts disulfures. Thèse de doctorat de l'Université Rennes 1, 2005.
- Ingrid Jacquemin & Jacques Nicolas. Modélisation de cystéines oxydées à l'aide de la programmation logique inductive. **JOBIM**, Lyon, juillet 2005, 331-340.
- Juliette Martin, Jean-François Gibrat & François Rodolphe. HMM for local protein structure. **ASMDA**, Brest, mai 2005, 180-187.
- Juliette Martin, Jean-François Gibrat & François Rodolphe. How to choose the optimal hidden Markov model for protein secondary structure prediction. **IEEE Intelligent Systems**, Special issue on Data Mining for Bioinformatics, accepté, à paraître en novembre/décembre 2005
- Khalid Benabdeselem, Christophe Geourjon, Yann Guermeur & Nicolas Sapay. Apprentissage automatique, application à la prédiction de la structure secondaire et tertiaire des protéines. Communication sur invitation présentée dans la session thématique : Bioinformatique II, **ASTI**, Clermont-Ferrand, octobre 2005.
- Khalid Benabdeselem, Gilbert Deléage & Christophe Geourjon. A neural network system based on structural alignment and clustering for proteins fold recognition. **ECCB**, Madrid, septembre 2005, 85-88.
- Khalid Benabdeselem, Gilbert Deléage & Christophe Geourjon. Cores extraction based neural network model for proteins fold recognition. **JOBIM**, Lyon, juillet 2005, 341-347.
- Yann Guermeur, A Lifchitz & Régis Vert. A kernel for protein secondary structure prediction. In **Kernel Methods in Computational Biology**, Editors : B. Schölkopf, K. Tsuda & Jean-Philippe Vert, The MIT Press, 2004, 193-206.