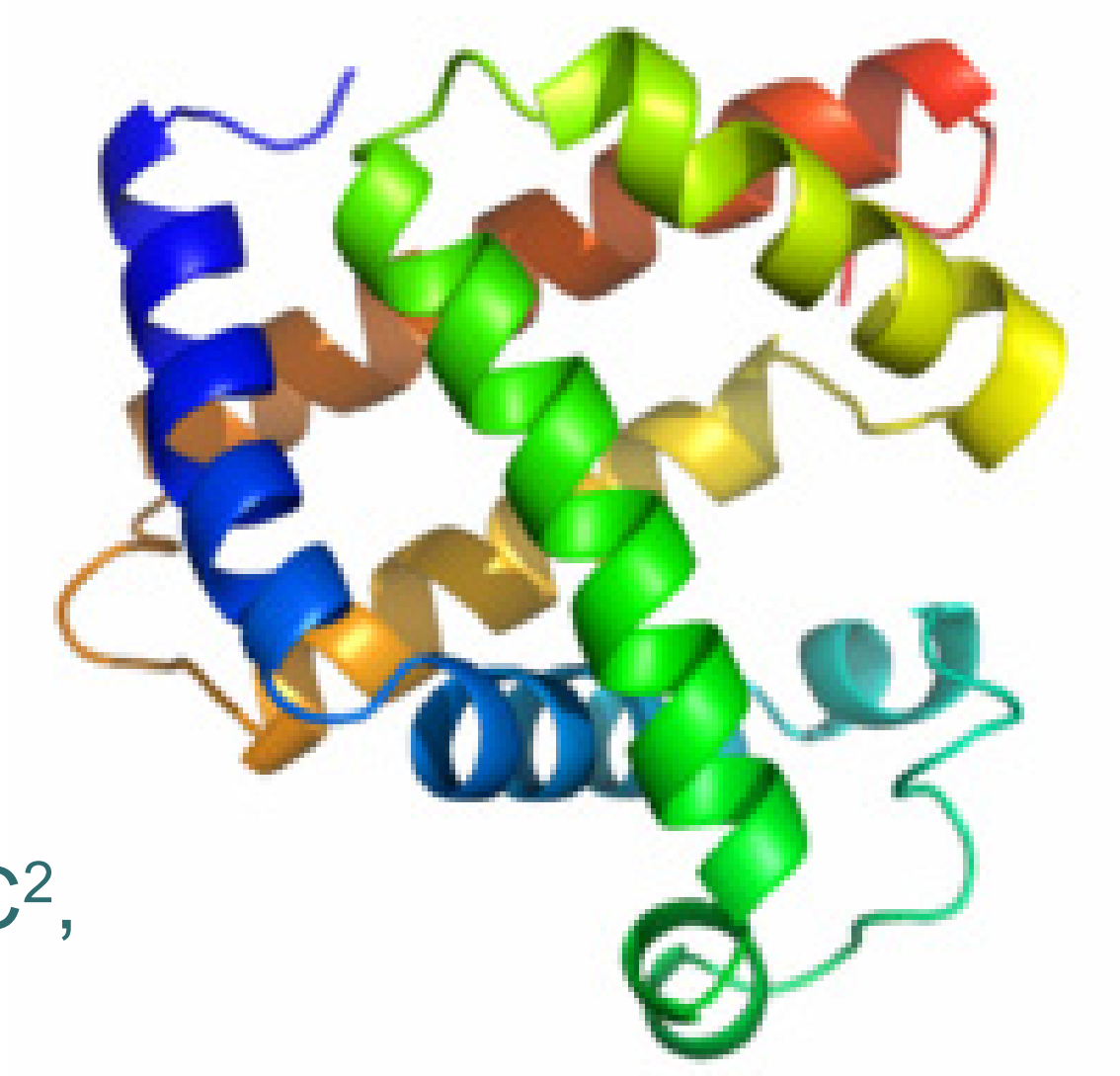


Projet GENOTO3D

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



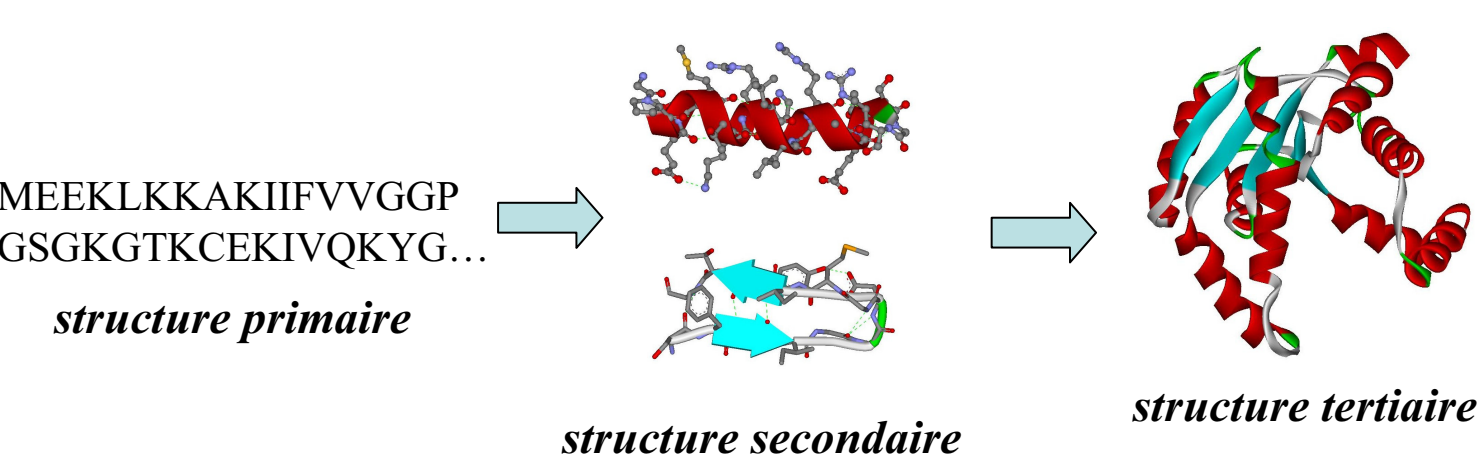
Guermeur Y¹, Benabdelsem K², Bréhélin L³, Capponi C⁴, Coste F⁶, Darcy Y¹, Deléage G², Denis F⁴, Gascuel O³, Geourjon C², Gibrat JF⁵, Jacquemin I⁶, Magnan C⁴, Marin A⁵, Martin J⁵, Monfrini E¹, Nicolas J⁶, Ralaivola L⁴, Taly JF⁵, Thomarat F¹

1. Equipe MODBIO (puis ABC) du LORIA (Nancy)
2. Laboratoire de Bioinformatique et RMN Structurales (LBRS) de l'IBCP (Lyon)
3. Equipe « Méthodes et algorithmes pour la bioinformatique » (MAB) du LIRMM (Montpellier)

4. Equipe « Bases de Données et Apprentissage Automatique » du LIF (Marseille)
5. Unité Mathématique, Informatique et Génome (MIG) de l'INRA, centre de Jouy-en-Josas
6. Projet Symbiose de l'IRISA (Rennes)

Contexte biologique :

Les **protéines**, macromolécules indispensables à la vie, assurent des fonctions très diverses. Cette **fonction biologique** est étroitement liée à la **structure 3D** de la protéine (cf. figure ci-dessous). La complexité de la détermination expérimentale de la structure 3D et la croissance exponentielle des données de séquences disponibles nécessitent la **mise au point de méthodes de prédiction de la structure 3D**. Ce problème central en biologie permet d'aborder l'essentiel des grandes questions ouvertes en traitement de données séquentielles.

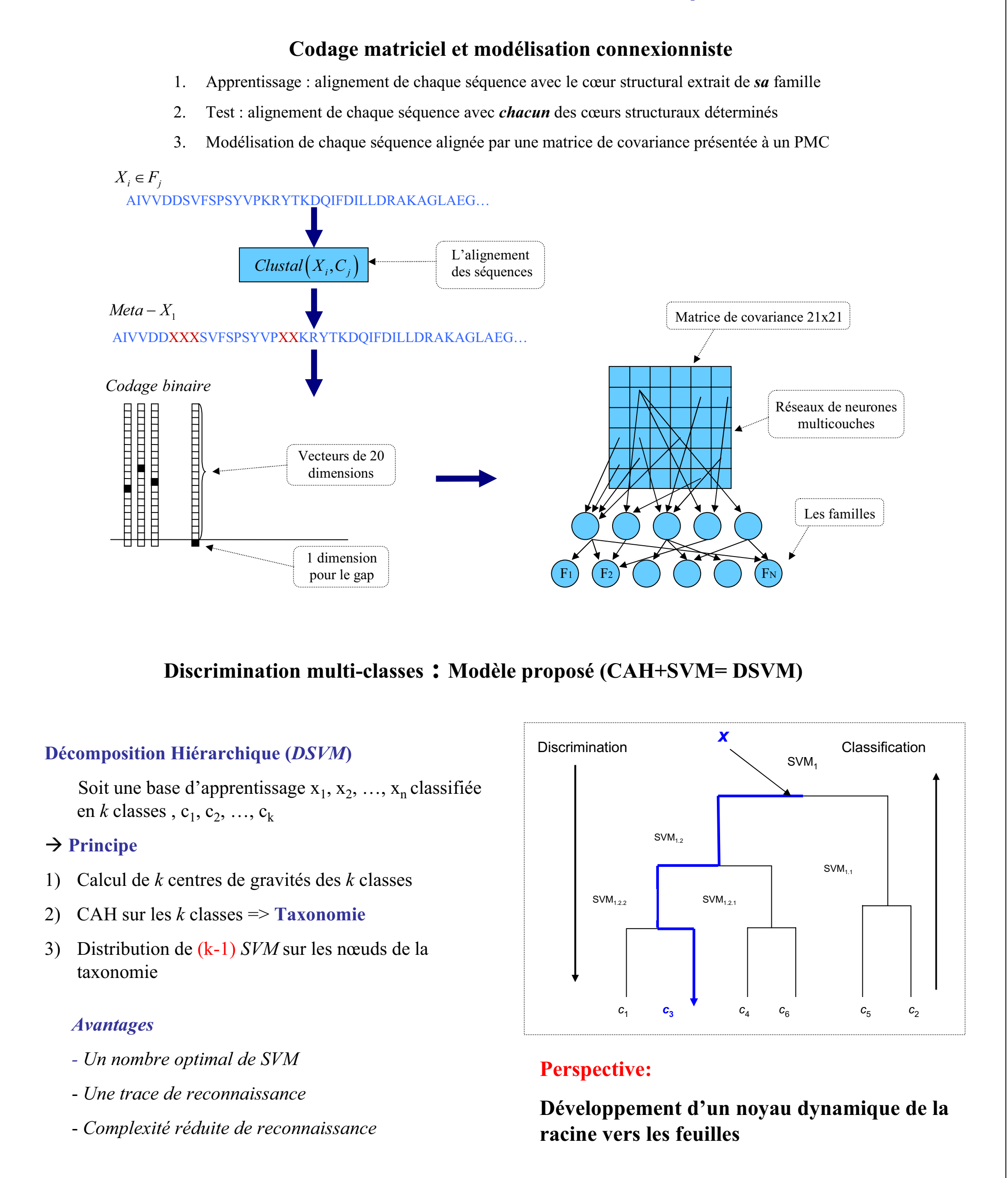
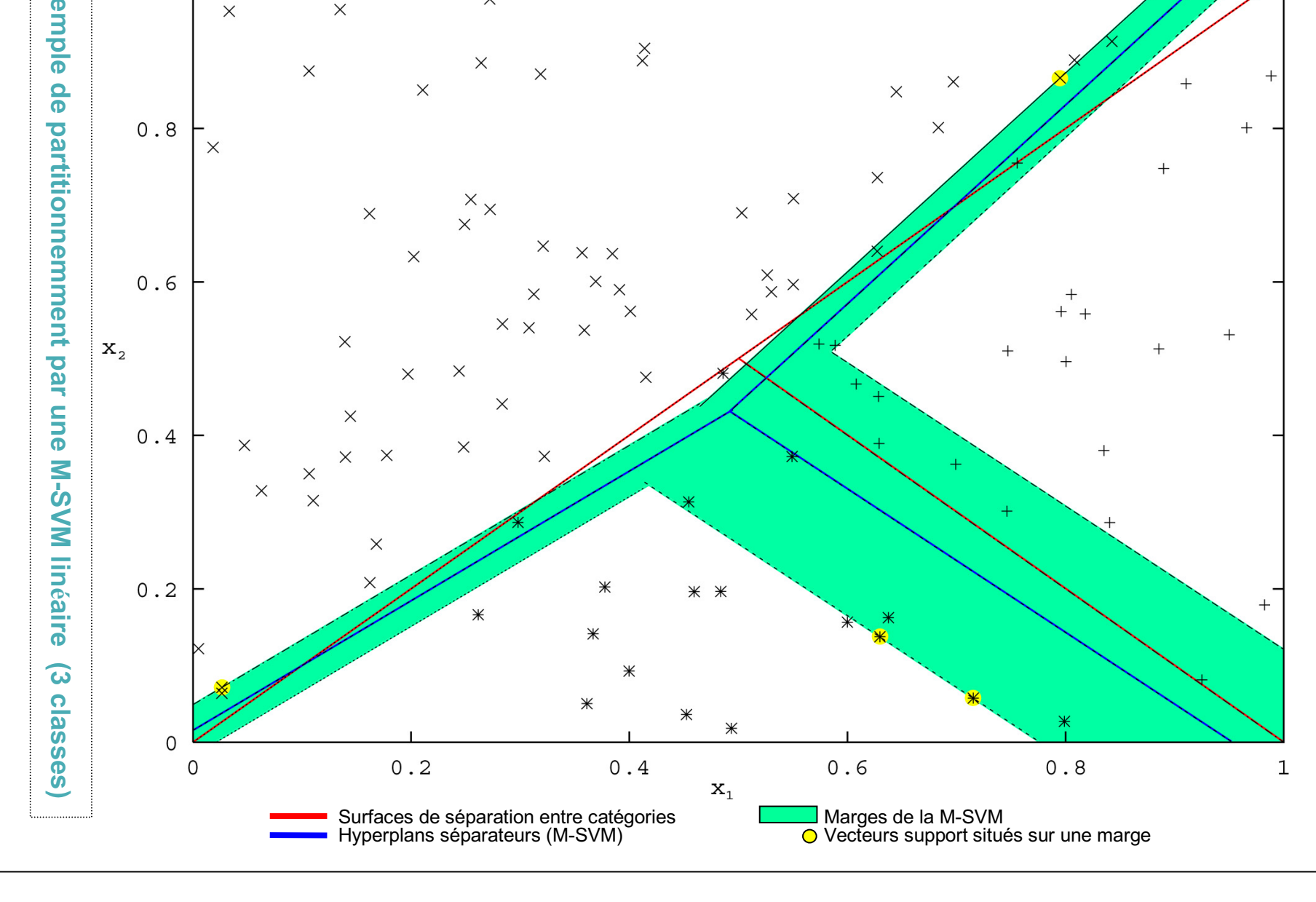
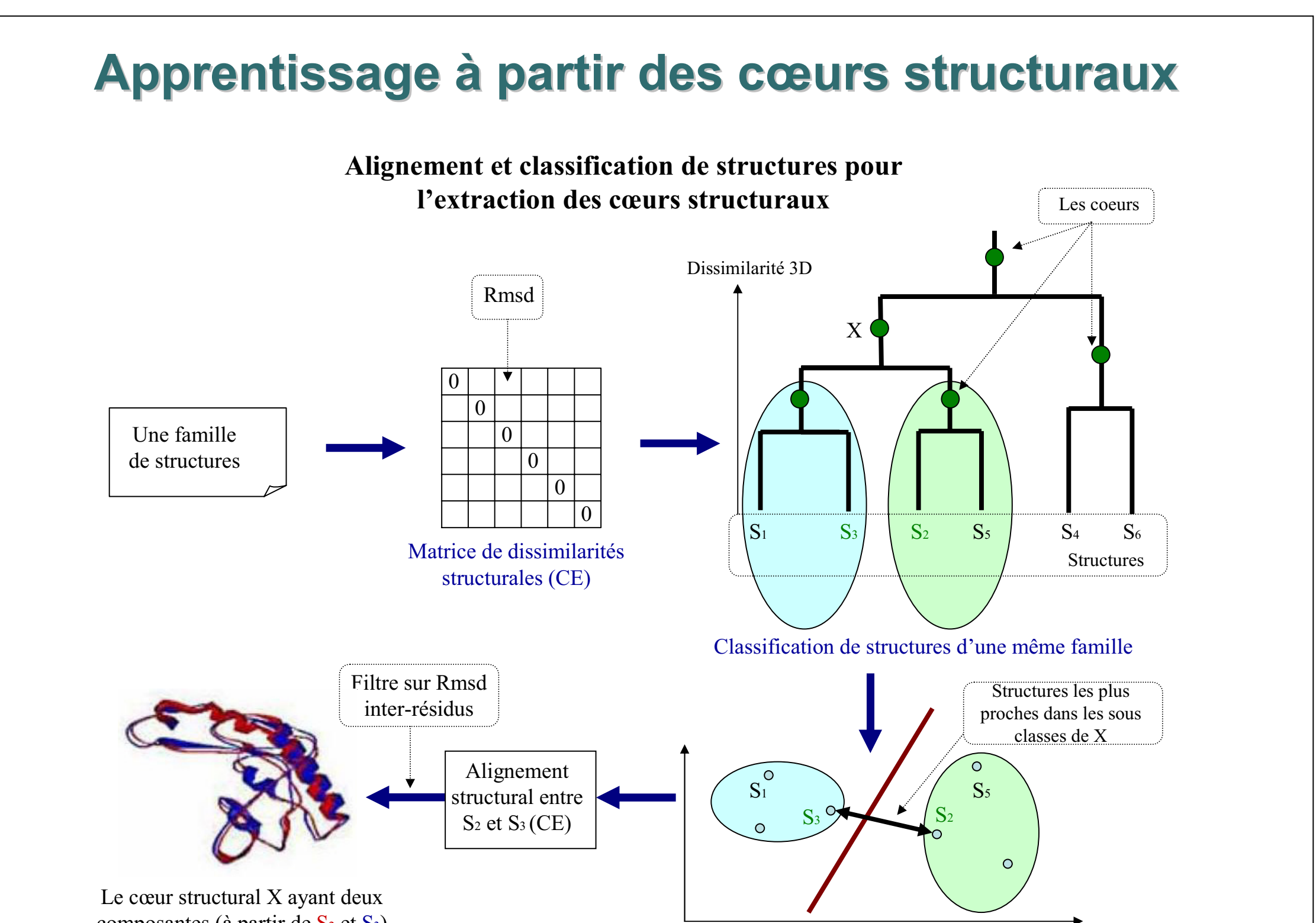
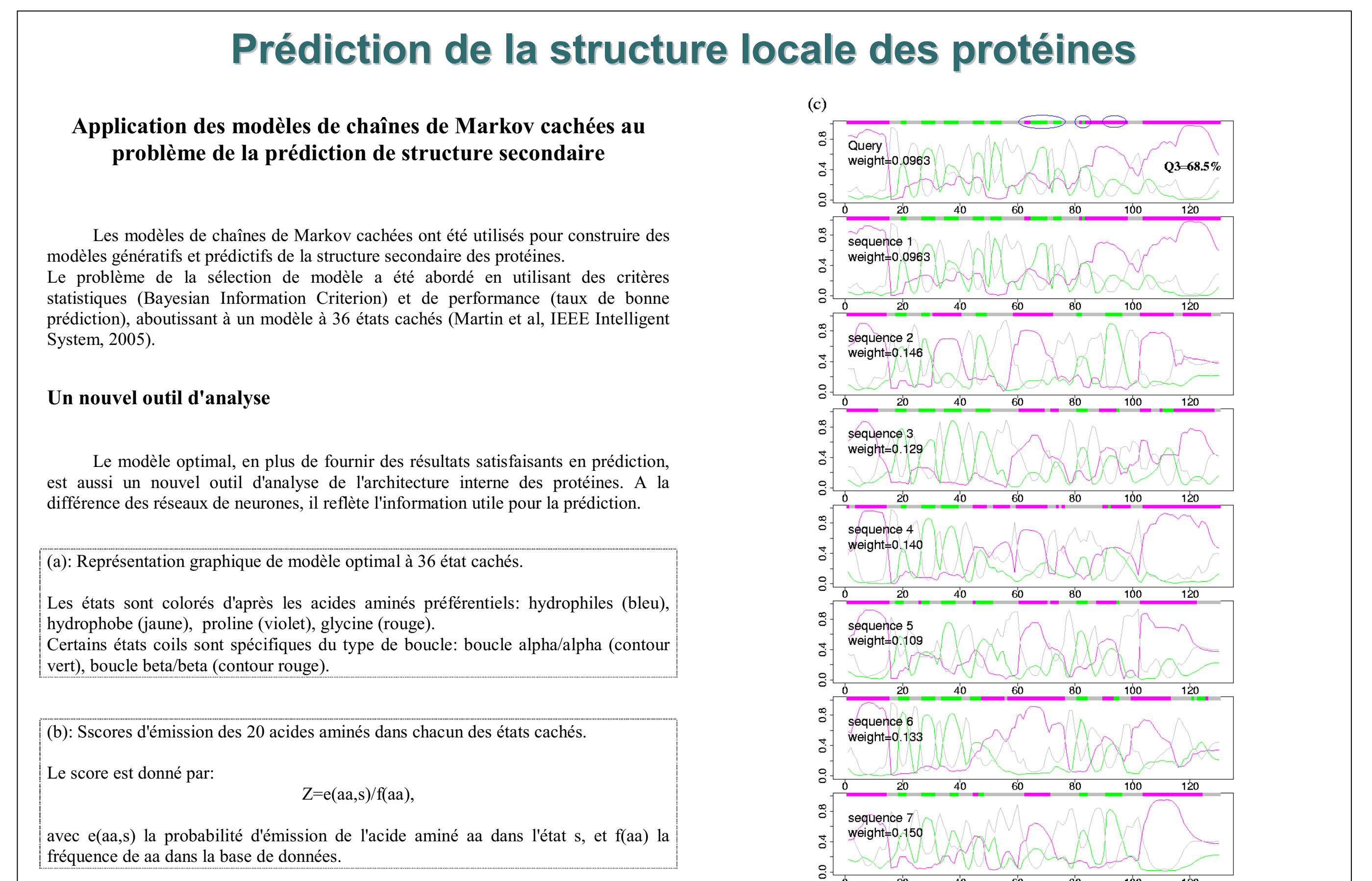
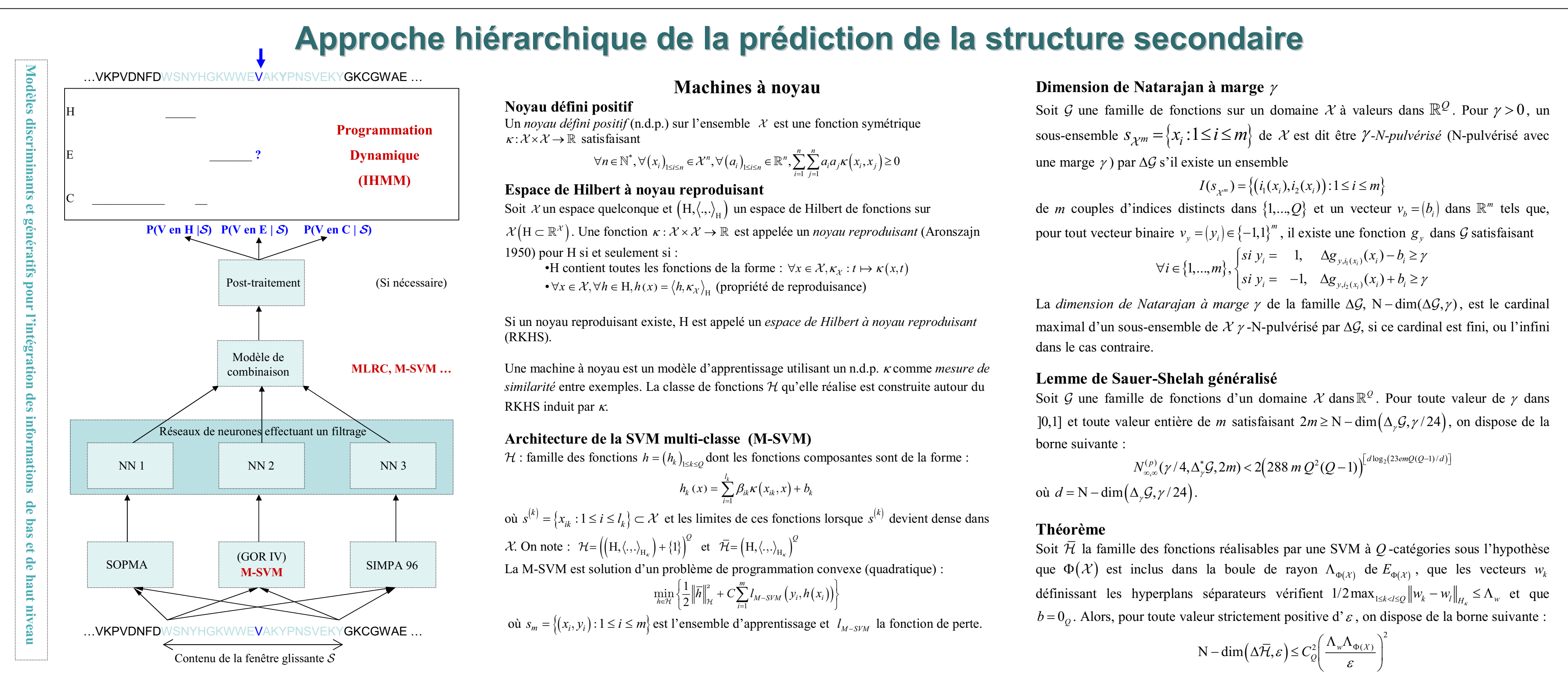


Les différents niveaux d'organisation structurelle des protéines : des régions de la **structure primaire** (ou séquence) de la protéine s'organisent en un ensemble d'éléments structuraux périodiques constituant la **structure secondaire** qui s'agencent dans l'espace pour former la **structure 3D** (*tridimensionnelle* ou *tertiaire*) de la protéine.

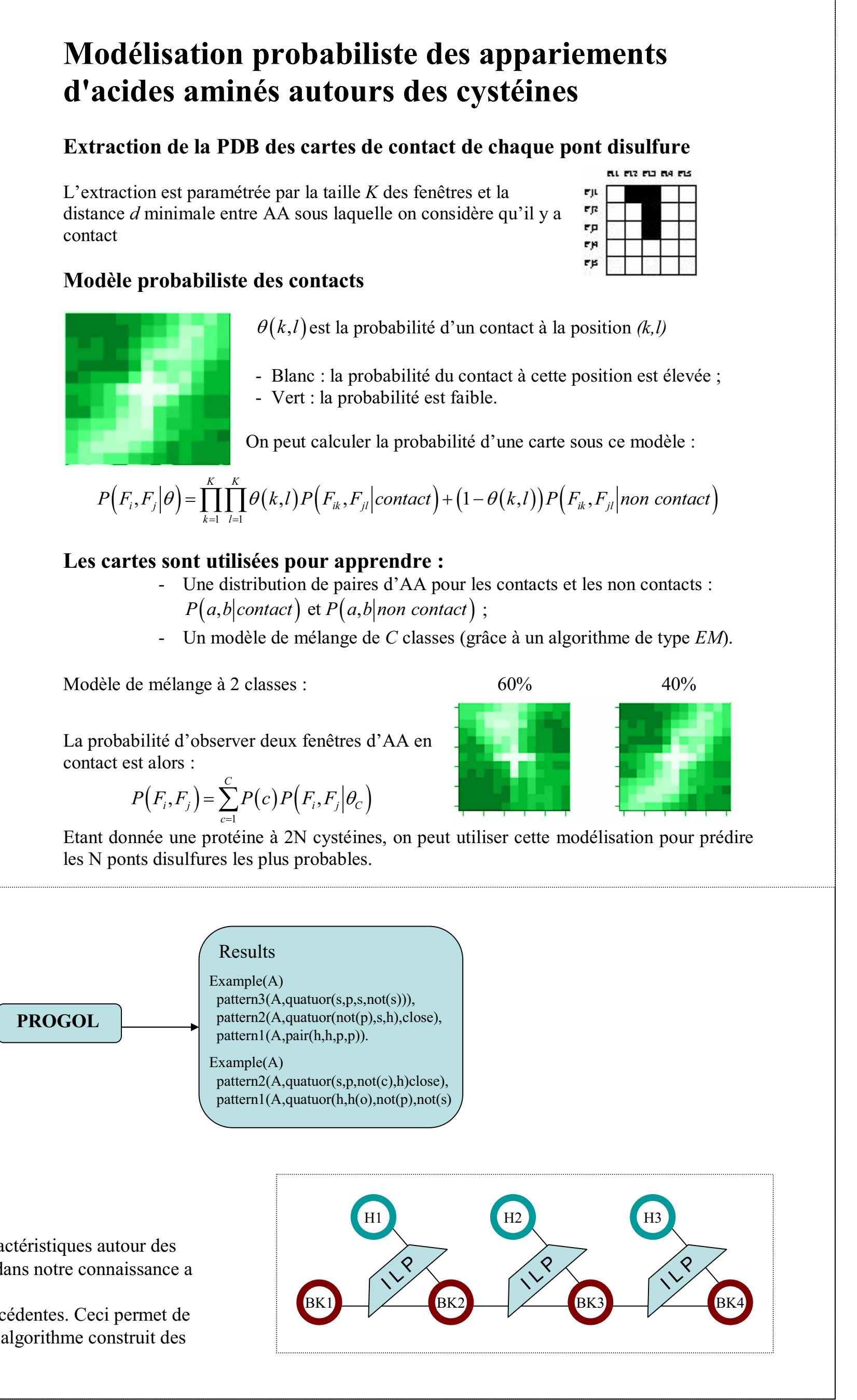
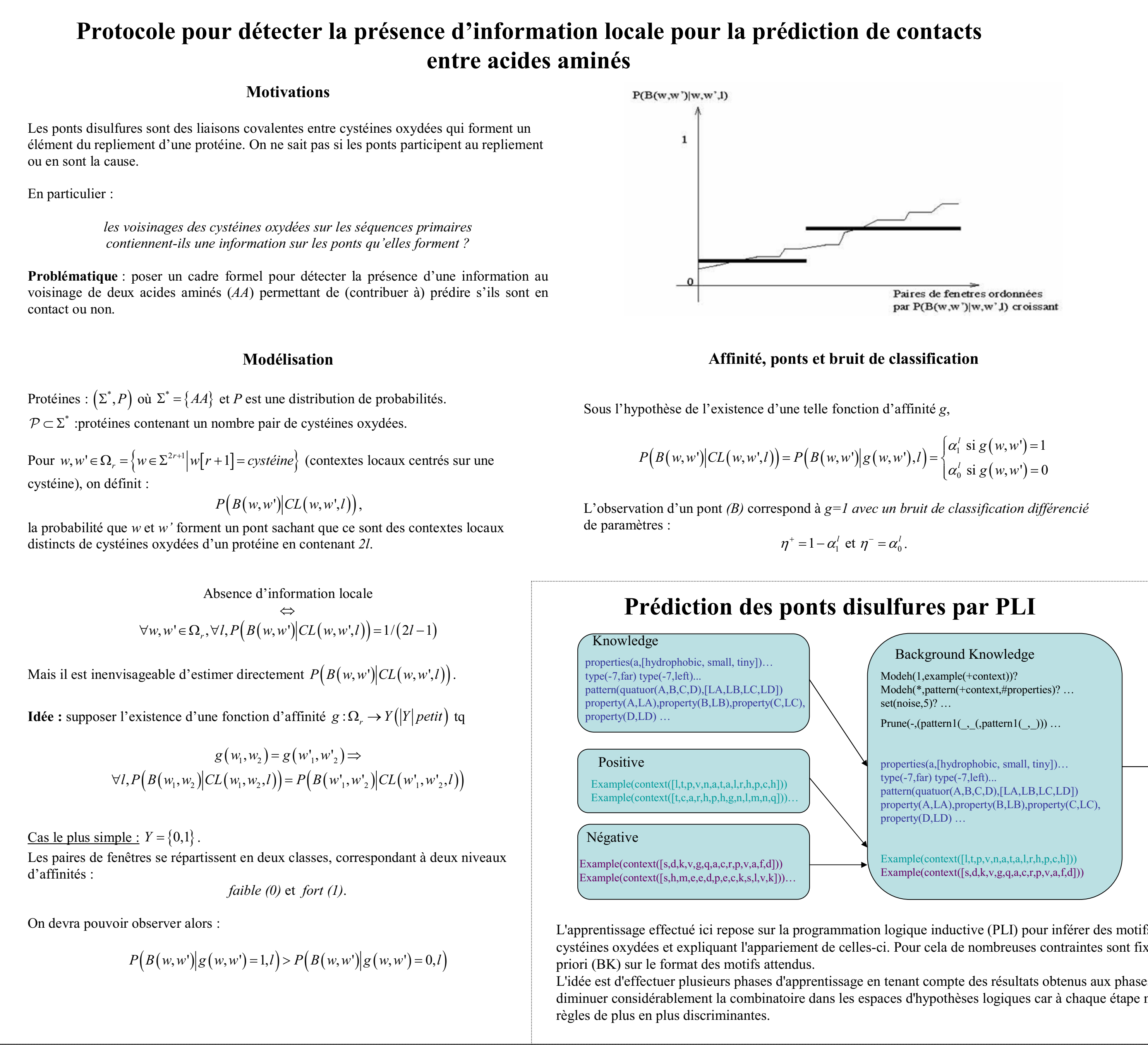
L'apprentissage automatique doit permettre d'exploiter les données de structures déjà existantes pour développer des méthodes de prédiction efficaces.

Le projet :

- **objectif** : dans le contexte de la prédiction de la structure tertiaire des protéines, mettre en évidence des problèmes de prédiction sur des séquences génériques et difficiles, et proposer des méthodes susceptibles de faire progresser l'état de l'art dans le domaine.
- **approche modulaire et hiérarchique** : ensemble de sous-problèmes et reformulation du problème
 - Prédiction des ponts disulfures et des ponts salins : IBCP, IRISA, LIF, LIRMM, LORIA
 - Prédiction de la **structure secondaire** (feuillettes β ...) : LORIA, MIG
 - Prédiction par **homologie** ou **analogie** et **reconnaissance des cœurs structuraux** : IBCP
 - Prédiction par **threading** : IRISA, MIG
 - Prédiction **ab initio** (*de novo*) : MIG



Prédiction des ponts disulfures



Références

- 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.
- Khalid Benabdeselem & Younés Bennani. Dendrogram based SVM for multi-class classification. *Journal of Computing and Information Technology - CIT 14*, 4, 2006, 291-297.
- Khalid Benabdeselem, Gilbert Deléage & Christophe Geourjon. Alignment structural et classification hiérarchique pour l'extraction des cœurs structuraux. Atelier : Extraction et gestion de connaissances appliquées aux données biologiques, *EGC-Lille*, Janvier 2006, 9-17.
- François Denis, Yann Esposito & Amary Habrad. Learning Rational stochastic languages. Proc. of the 19th Annual Conference on Learning Theory, *LNAI*, 4005, 2006, 274-288.
- Elodie Duprat, Marie-Paule Lefranc & Olivier Gascuel. A simple method to predict protein-binding from aligned sequences - application to MHC superfamily and 2-microglobulin. *Bioinformatics*, 2006, 22, 4, 453-459.
- Yann Guermeur. Large Margin Multi-category Discriminant Models and Scale-sensitive Psi-dimensions. *Rapport de recherche INRIA*, RR-5314, 2004 (révisé en 2006).
- 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.
- Christophe Magnan. Asymmetrical Semi-Supervised Learning and Prediction of Disulfide Connectivity in Proteins. *RIA*, 2006, à paraître.
- Juliette Martin, Jean-François Gibrat & François Rodolphe. Analysis of an optimal hidden Markov model for secondary-structure prediction. *Soimis*.
- Juliette Martin, Jean-François Gibrat & François Rodolphe. Choosing the optimal hidden Markov model for secondary-structure prediction. *IEEE Intelligent Systems*, 20, 2005, 19-25.
- Liva Ralaivola, François Denis & Christophe Magnan. CN-CPN. Proc. of the 23rd Int. Conf. on Machine Learning, 2006, 721-728.
- Nicolas Sapay, Yann Guermeur & Gilbert Deléage. Prediction of amphipathic in-place membrane anchors in monogenic proteins using a SVM classifier. *BMC Bioinformatics*, 2006, Vol.7, 255.
- Raluca Uricaru, Eric Rivais & Laurent Bréhélin. Hidden Markov models for the detection of motifs repeats in protein sequences. *IPG*, Lyon, novembre-décembre 2006.

