

Modélisation Statistique de l'Environnement : Réduction par Classifications

Charantonis Anastase Alexandre

Post-doc sous Fouad Badran

Equipe Méthodes Statistiques de Data
Mining et Apprentissage

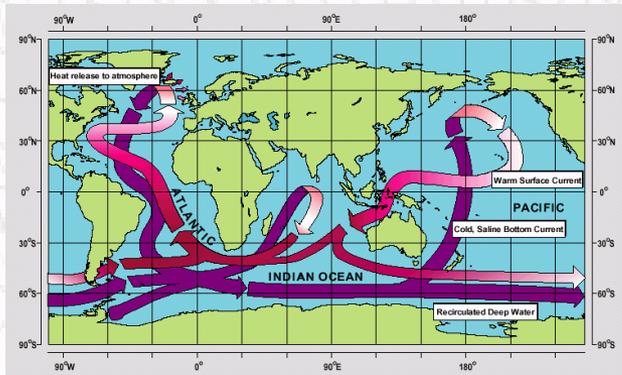
Paris, 11-6-2015

- 
- **Introduction**
 - *La Méthodologie Statistique*
 - *Applications*
 - *Complétion de données*
 - *Conclusions - Perspectives*

INTRODUCTION

Modélisation de l'environnement

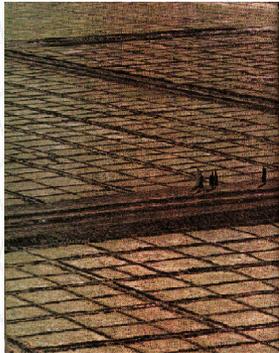
L'environnement est un *systeme complexe*. Essayer de bien comprendre les choses est un réflexe induit par la curiosité humaine, mais bien comprendre notre environnement est crucial car cela permet *d'anticiper* son évolution future.



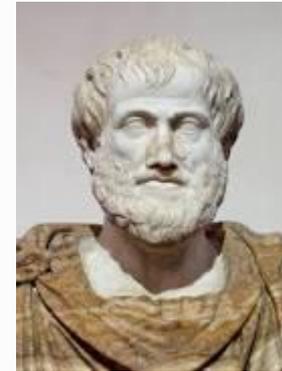
Une meilleure compréhension induit une meilleure capacité de prédire les répercussions de nos actions, et peut-amener une amélioration du niveau de vie.

INTRODUCTION

Modélisation de l'environnement



Egyptian field, from Ancient Egypt, page 48 (see references).



Etudes de l'environnement existent depuis l'antiquité. Nous n'avons commencé à archiver de façon précise des données environnementales que depuis "récemment" (200 ans de données de qualité croissante).

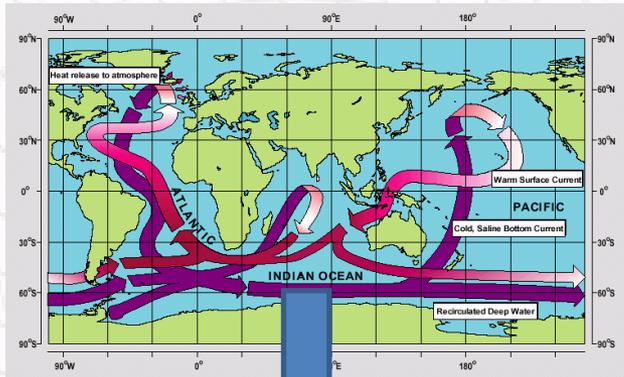


Que faire de ces données?

INTRODUCTION

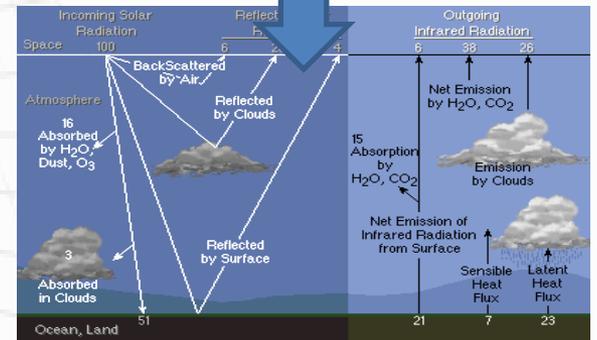
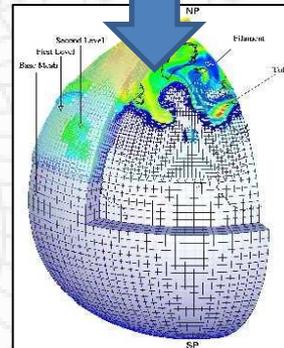
Modélisation de l'environnement

Premier instinct: identifier les lois du système.



$$\rho \left(\frac{\partial \mathbf{v}}{\partial t} + \mathbf{v} \cdot \nabla \mathbf{v} \right) = -\nabla p + \nabla \cdot \mathbf{T} + \mathbf{f}$$

(Navier-Stokes)



Pourtant, la **haute complexité** et la **non-linéarité** du système impliquent que, bien qu'on puisse le reproduire en se basant sur des systèmes d'équations différentielles, de **faibles variations** des valeurs initiales peuvent induire des évolutions **fortement divergentes**.

INTRODUCTION

Modélisation de l'environnement

Modélisation numérique de l'environnement:

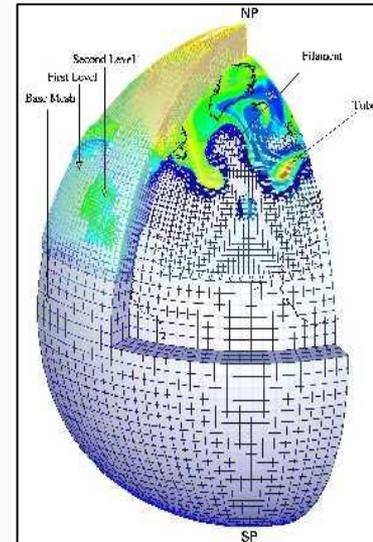
Nécessite:

- Connaissance explicite des équations qui régissent son fonctionnement.
- Ressources informatiques fortes.
- Paramétrisation et initialisation très fines du modèle.

Fournit:

- Prévision de l'évolution temporelle du système.

Faibles variations des valeurs initiales induisent des évolutions divergentes.

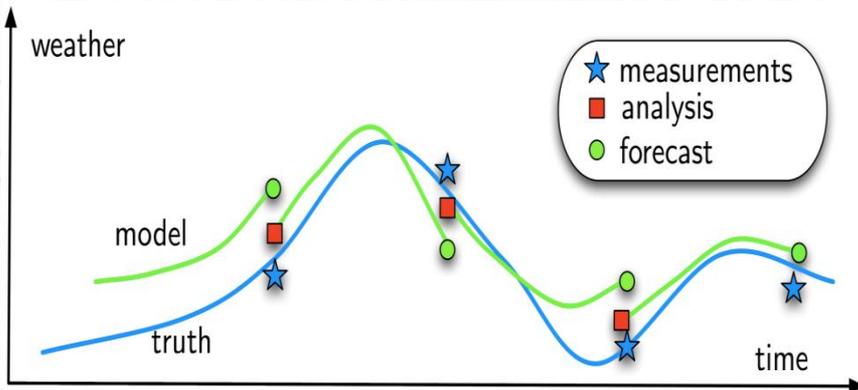
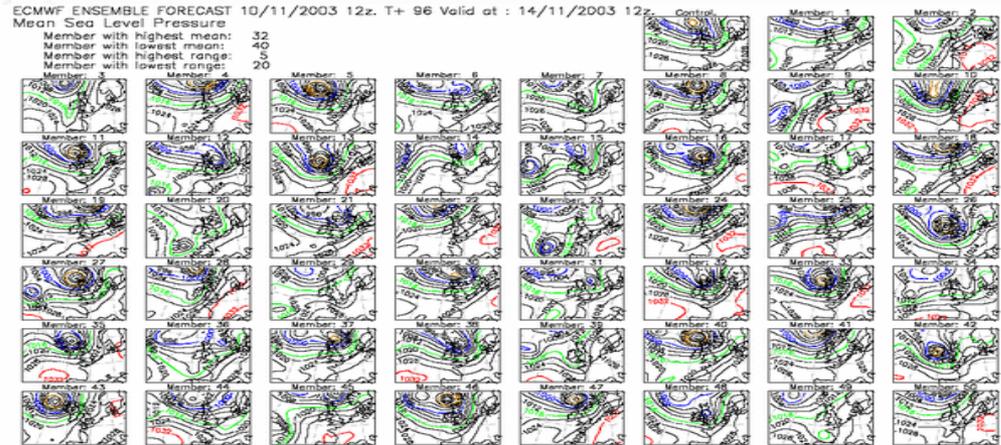


INTRODUCTION

Modélisation de l'environnement

Prévision d'ensembles:

Des conditions initiales légèrement perturbées génèrent un groupe d'évolutions possibles du système.

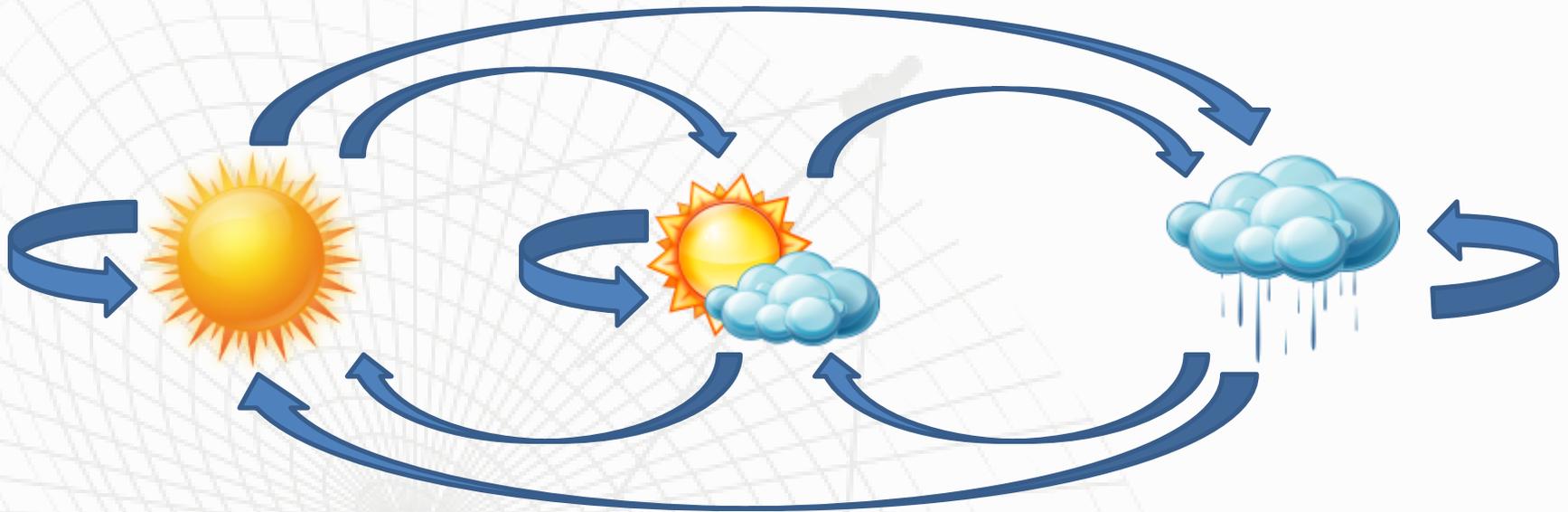


Assimilation des données:

Modification des paramètres et/ou des conditions initiales d'un modèle numérique pour mieux coller aux observations du système, tout en prenant en compte les incertitudes du modèle et des données.

INTRODUCTION

Classifier l'environnement



Si on a suffisamment de données, on peut réduire un système environnemental en un *ensemble de classes*, chacune représentant une des *situations* possible du système. Chaque situation a une probabilité *d'évoluer* dans le temps en un autre situation. Nous pouvons estimer ces probabilités à partir de séquences de données.

Compromis nécessaire entre la finesse des classes et la qualité des probabilités calculées.

INTRODUCTION

Modélisation de l'environnement

PROBLEMES CONCRETS



MODELISATION
STATISTIQUE



GENERALISATION DE
LA METHODE

OCÉANOGRAPHIE



PROFHMM

(Profile reconstruction through
Hidden Markov Models)

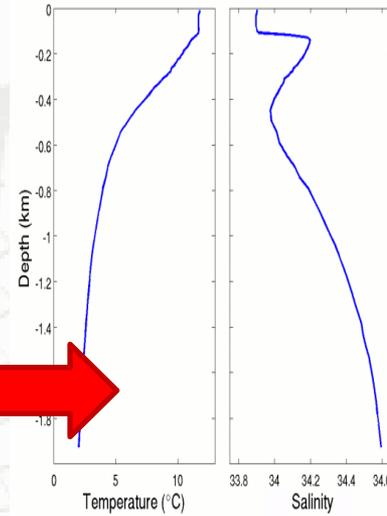
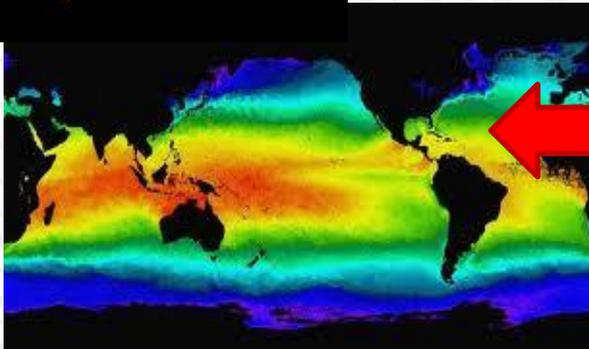


AUTRES APPLICATIONS

(Météorologie, Climatologie, Hydrologie
Comportement Sociaux...)

INTRODUCTION

Les observations océaniques

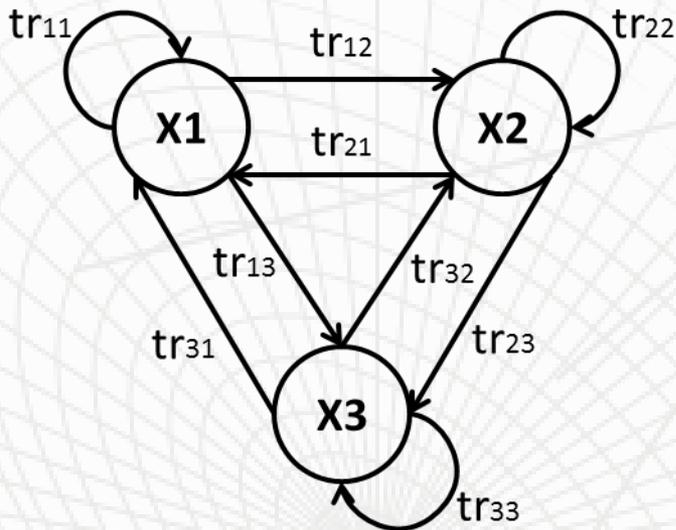


OBSERVATIONS SATELLITE:
Grande couverture spatiale
Grand nombre de variables
Haute fréquence d'images
Observation d'émissions de surface

PROFONDEURS:
Mise en place de mesures in-situ
(bouées, gliders, flotteurs,
campagnes océanographiques):
Irrégulières en temps et espace
Mise en place difficile.

- *Introduction*
- ***La Méthodologie Statistique***
 - Modèles de Markov Cachés (HMM)
 - Cartes Topologiques Auto-Organisatrices
 - Algorithme de Viterbi
- *Applications*
- *Complétion de données*
- *Conclusions - Perspectives*

MODELES DE MARKOV



Chaîne de Markov
à 3 états

Chaque état:

Probabilité définie de passer à chaque état du système à chaque pas de temps.

Propriété de Markov de premier ordre:

$$P(X_{i_t}^t | X_{i_{t-1}}^{t-1} X_{i_{t-2}}^{t-2} \dots X_{i_0}^0) = P(X_{i_t}^t | X_{i_{t-1}}^{t-1})$$

Et elle est égale à:

$$tr_{i_t, i_{t-1}} = P(X_{i_t}^t | X_{i_{t-1}}^{t-1})$$

(Probabilité de transition)

Connaissance, à tout moment, de l'état du système.

MARKOV: HYPOTHESES

Etats discrets:

On considère un système océanique simplifié représenté par un *ensemble discret de profils verticaux*.

Propriété de Markov de premier ordre:

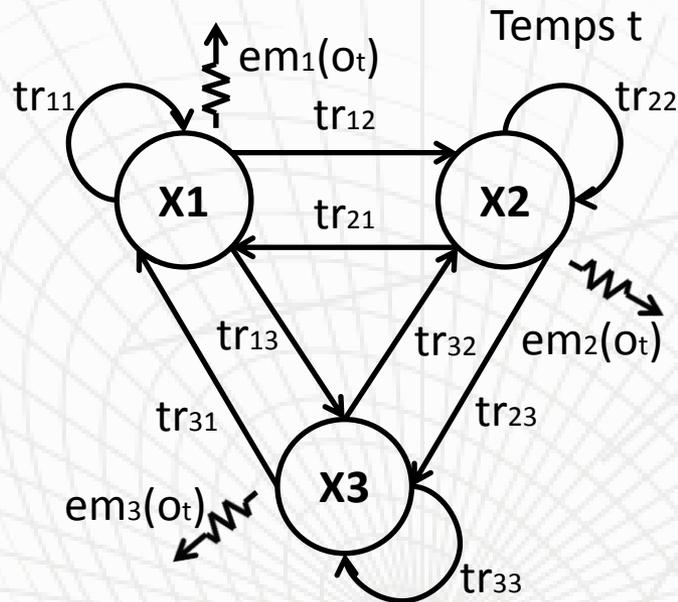
Etat de l'océan à **$t+1$** dépend seulement de l'état à **t** .

Les probabilités sont homogènes dans le temps:

Evolution pluriannuelle lente.

- QUESTIONS:**
1. **Comment retrouver les profils verticaux à partir des observations de surface?**
 2. **Comment déterminer les états?**

MODELES DE MARKOV CACHES



Chaîne de Markov
Cachée à 3 états

Probabilités initiales: $\pi_i = P(X_i)$

Etats du système non-observable directement:
Modèle de Markov caché
(Hidden Markov Model)

A chaque pas de temps :

- **une observation liée au système,**
- **sa probabilité d'être émise par chacun des états du système.**

$$em_{j_t}(o_t) = P(o_t | X_{j_t}^t)$$

(Probabilité d'émission)

États, em , tr , π & série d'observations

Série d'états *cachés* la plus probable d'avoir
générée ces observations

2. Comment générer les états?

CARTES TOPOLOGIQUES AUTO-ORGANISATRICES

Méthode de Classification
Statistique Neuronale

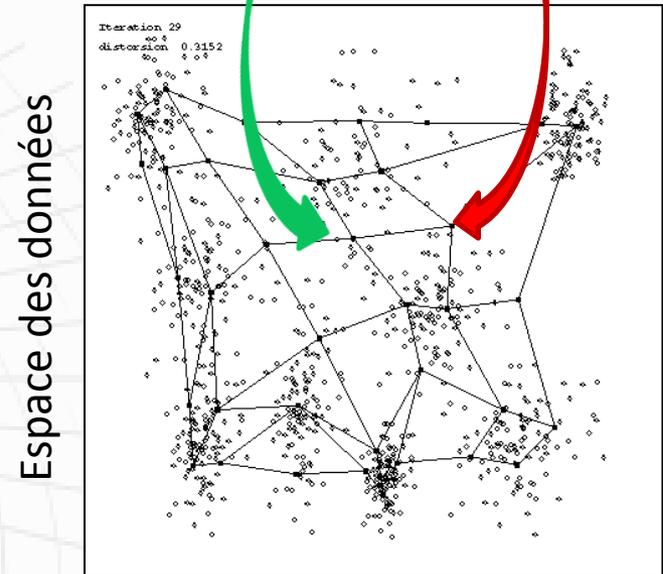
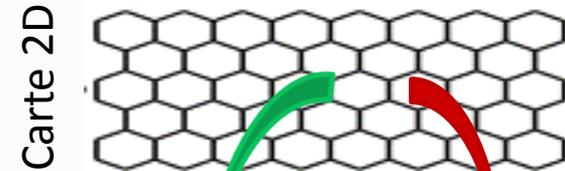
Entrées: Données Multidimensionnelles.

Sortie: *Projection de ces données de dimension supérieure sur une carte (généralement) 2D.*

Regroupements en **classes** représentés par:

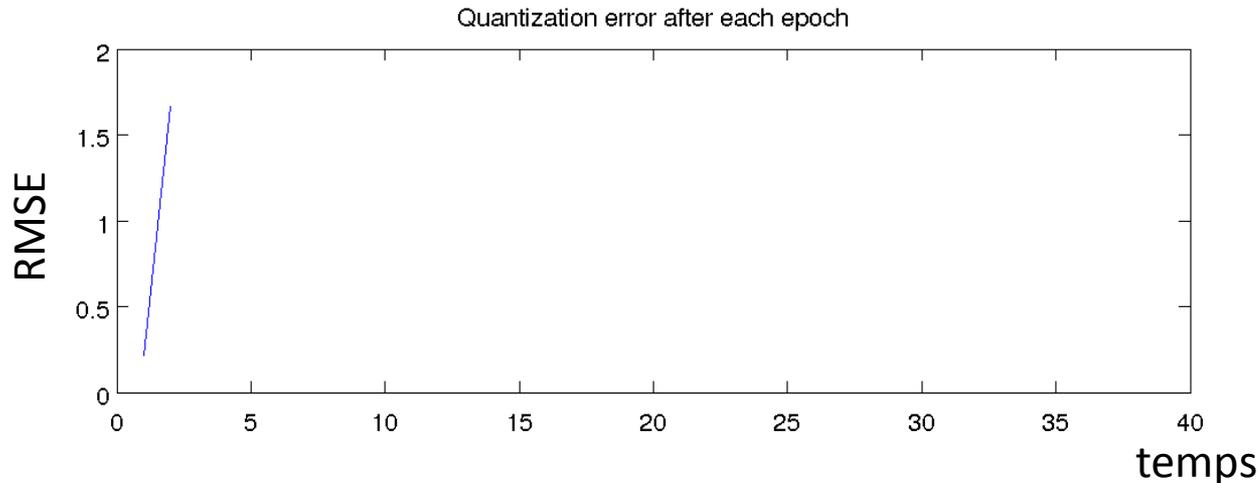
- un **réfèrent** dans l'espace des données
- un **index** sur la carte 2D

La carte est déterminée après apprentissage à partir de données.

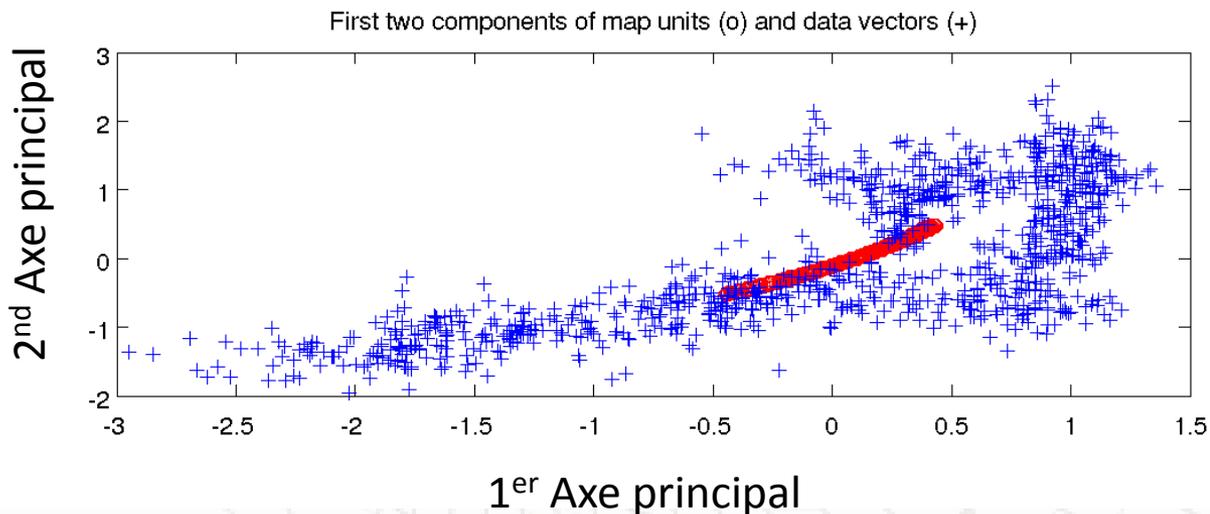


Deux classes voisines sur la carte 2D  Deux situations proches dans l'espace des données

CARTES TOPOLOGIQUES: ENTRAINEMENT



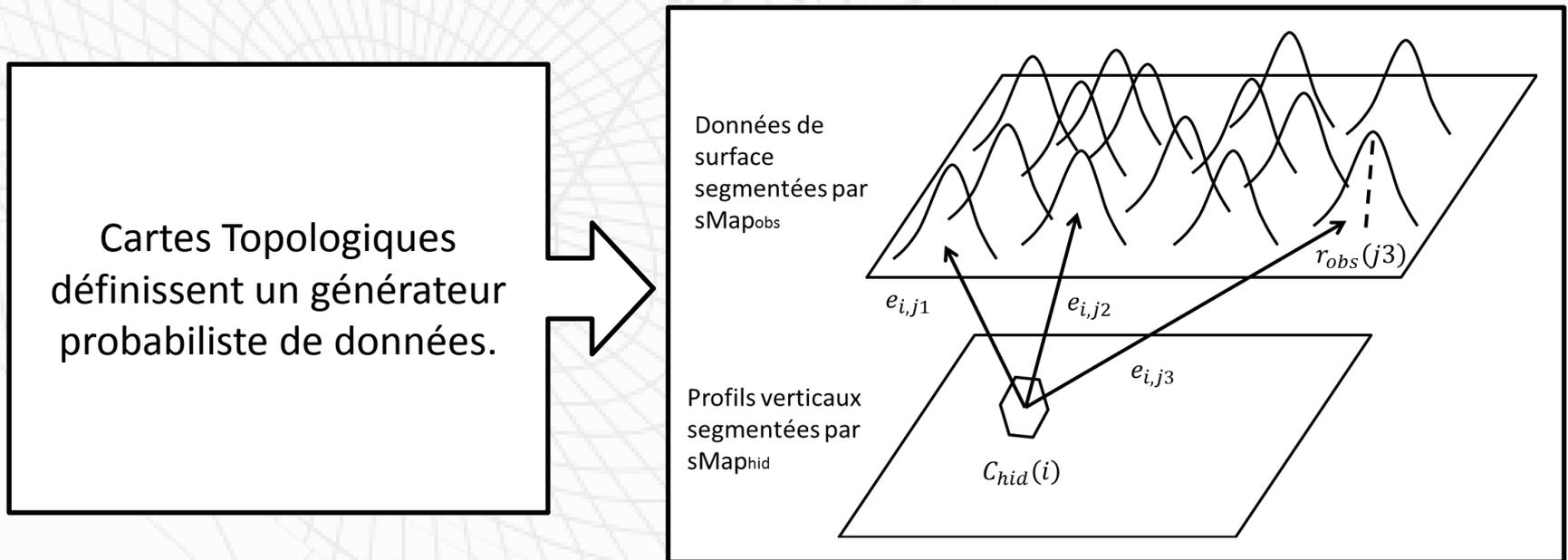
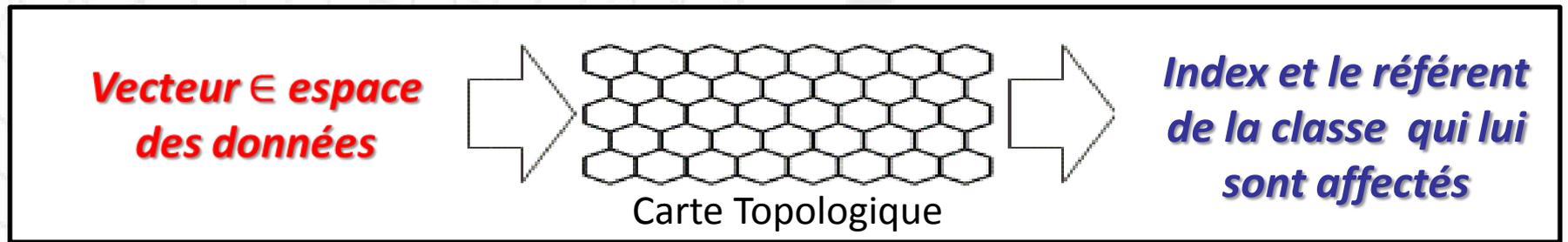
Entraînement d'une carte topologique avec 600 états-types, à partir de 1065 vecteurs de dimension 5.



Projections, durant l'entraînement des données et des états types sur le premier plan principal de l'ACP.

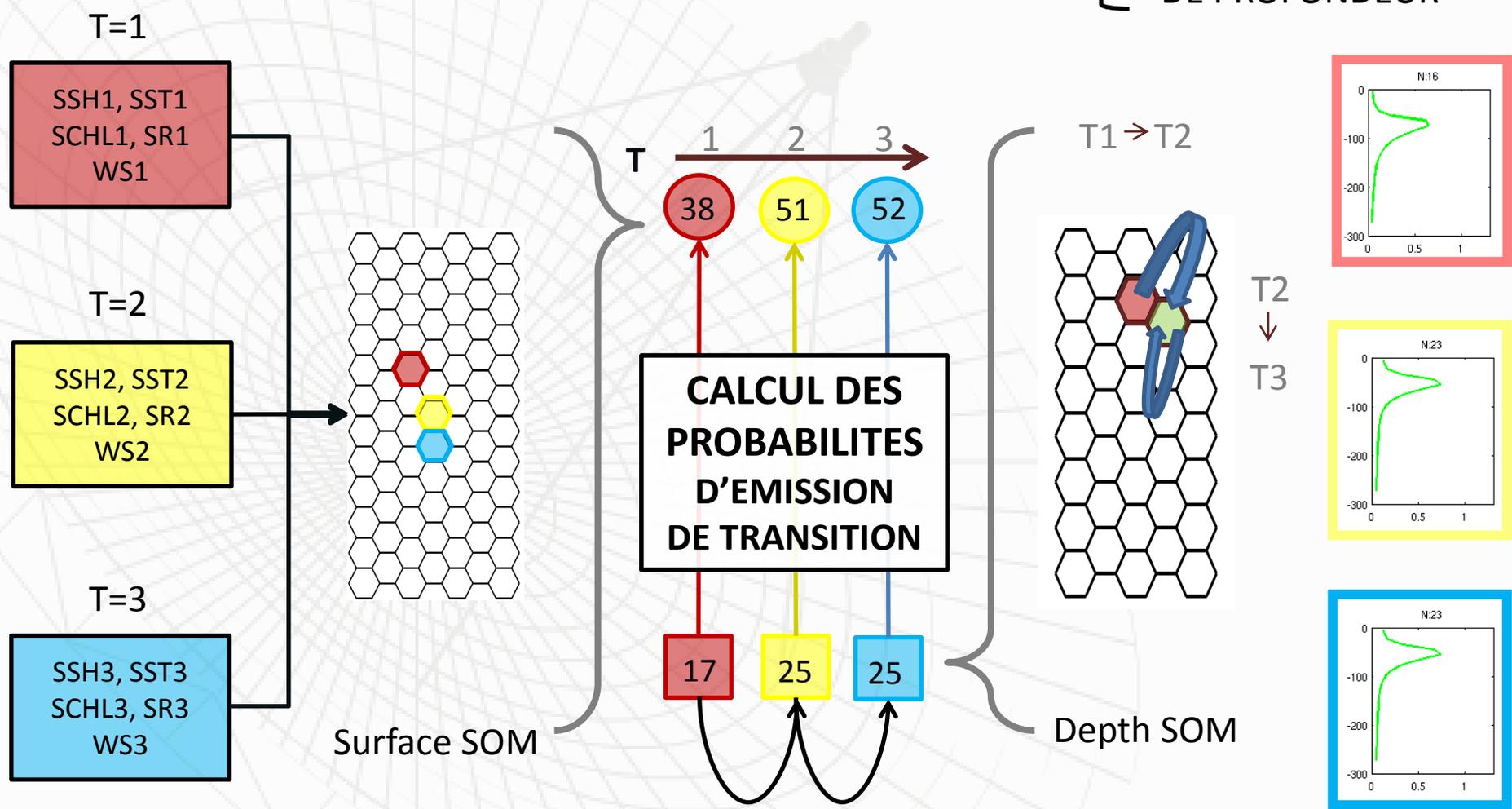
CARTES TOPOLOGIQUES: CLASSIFICATION ET GENERATEURS

CLASSIFICATION - DISCRETISATION



PROFHMM: REDUCTION DU PROBLEME MULTIDIMENSIONNEL

Nous obtenons des SERIES D'INDEX } DE SURFACE
DE PROFONDEUR



PROFHMM: UTILISATION DES PROPRIETES TOPOLOGIQUES

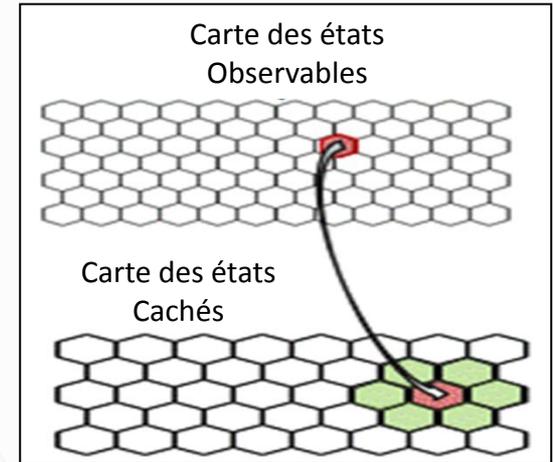
UTILISATION DE LA NOTION DE **VOISINAGE** DES CARTES TOPOLOGIQUES *POUR AMELIORER LES* PROBABILITES DES HMM

$$Em_{final}(i, j) = \sum_{N_{seq}} \left(\frac{w_c * \sqrt{L_{seq}} * Em_{B-W_{seq}}(i, j) + \sum_{k=1}^{N_{hid}} (NM_{hid}(j, k) * Em_{B-W_{seq}}(i, k))}{\sum_{k=1}^{N_{hid}} (NM_{hid}(j, k) * Em_{B-W_{seq}}(i, k))} \right) + 1$$

Normalisé pour $\sum_{i=1}^{N_{hid}} e_{i,j} = 1$

$$Tr_{final}(i, j) = \sum_{N_{seq}} \left(\frac{w_c * \sqrt{L_{seq}} * Tr_{B-W_{seq}}(i, j) + \sum_{k=1}^{N_{obs}} (NM_{hid}(i, k) * Tr_{B-W_{seq}}(i, k))}{\sum_{k=1}^{N_{obs}} (NM_{hid}(i, k) * Tr_{B-W_{seq}}(i, k))} \right) + 1$$

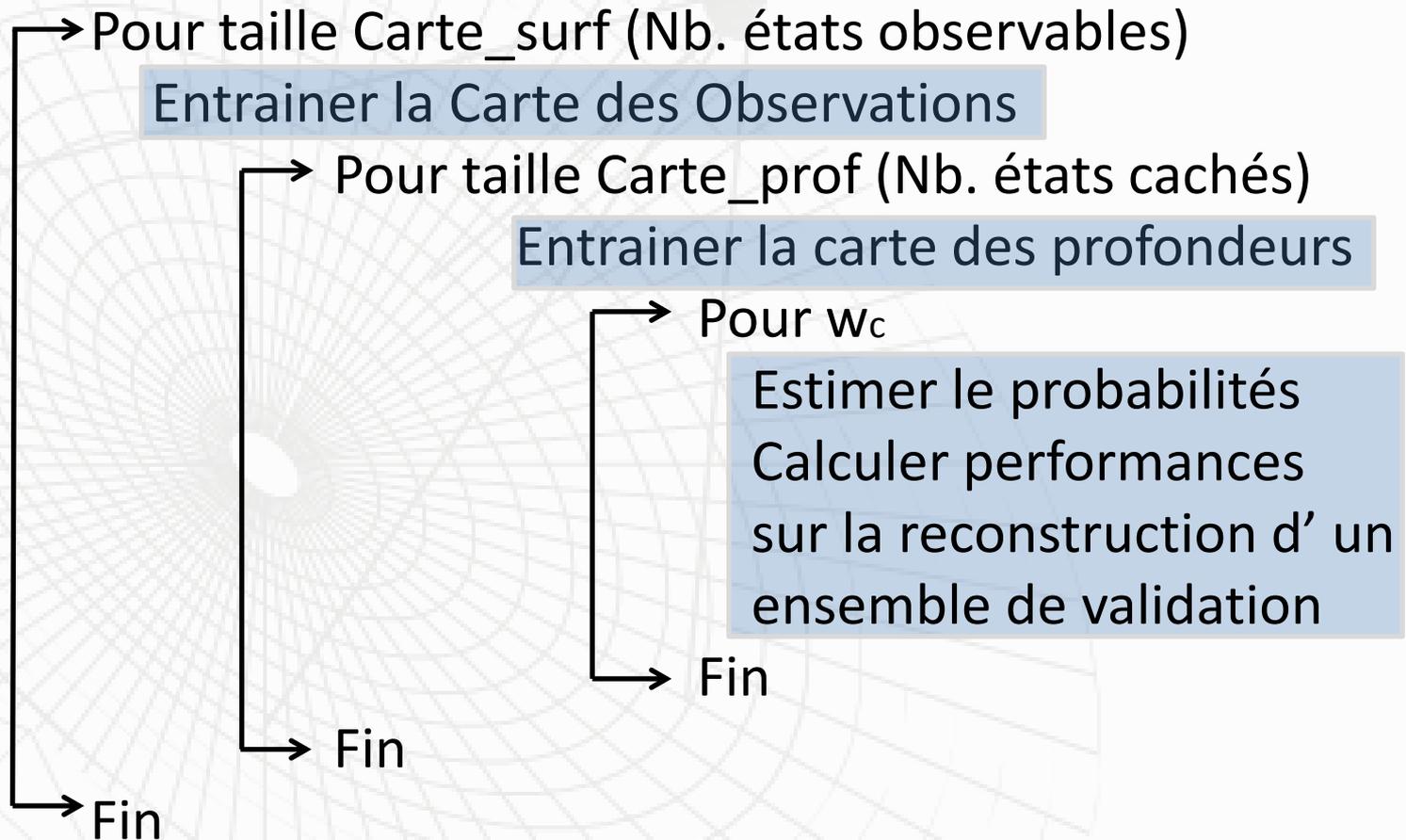
Normalisé pour $\sum_{i=1}^{N_{hid}} tr_{i,j} = 1$



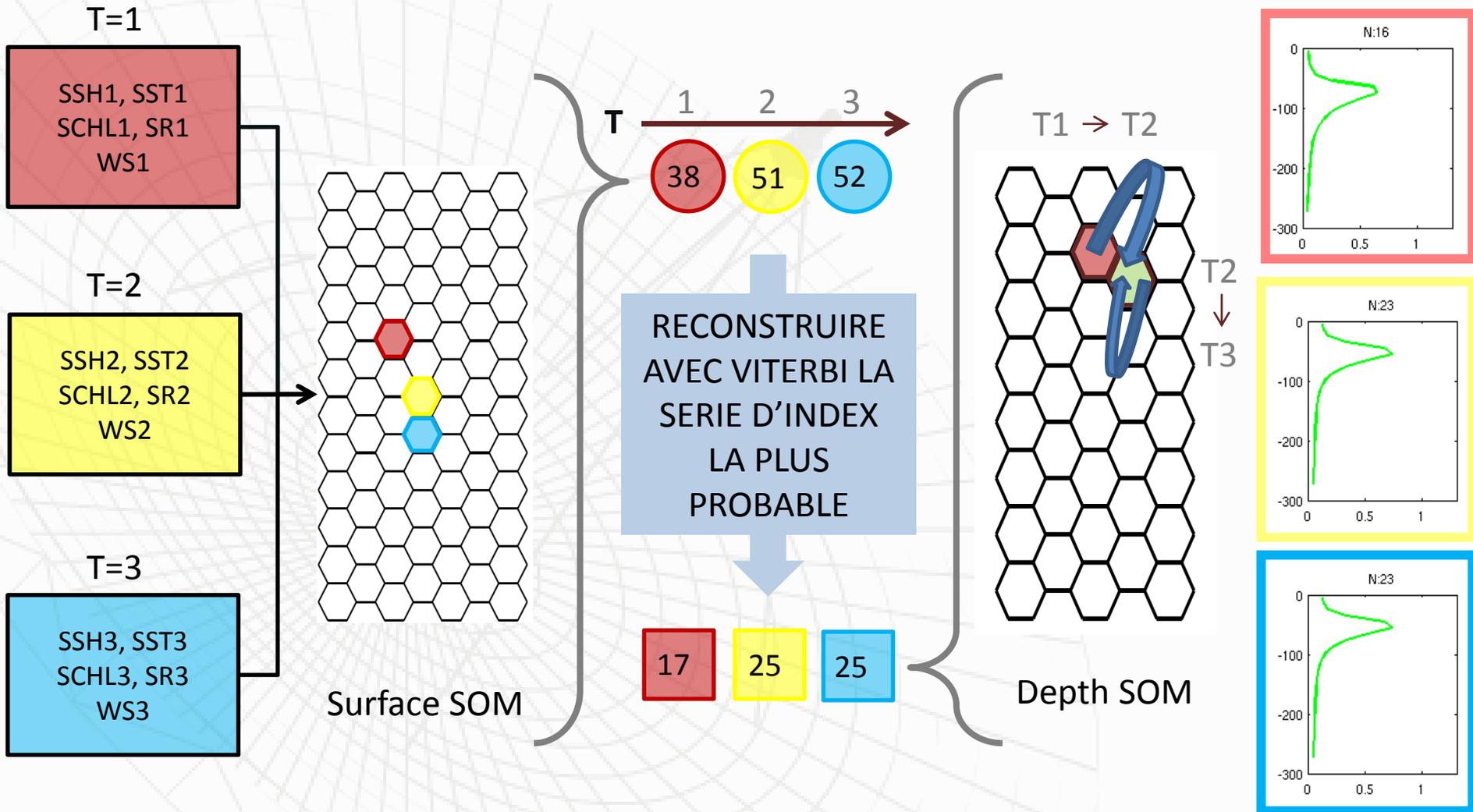
N_{hid}	Nombre d'états cachés
w_c	Terme de pondération
N_{obs}	Nombre d'états observables
N_{seq}	Nombre de sequences d'entraînement
L_{seq}	Longueur de chaque sequence utilisée
$Tr_{B-W_{seq}}$	Matrices calculées sur chaque séquence
$Em_{B-W_{seq}}$	(compteurs/Baum-Welch)

PROFHMM: MISE EN ŒUVRE

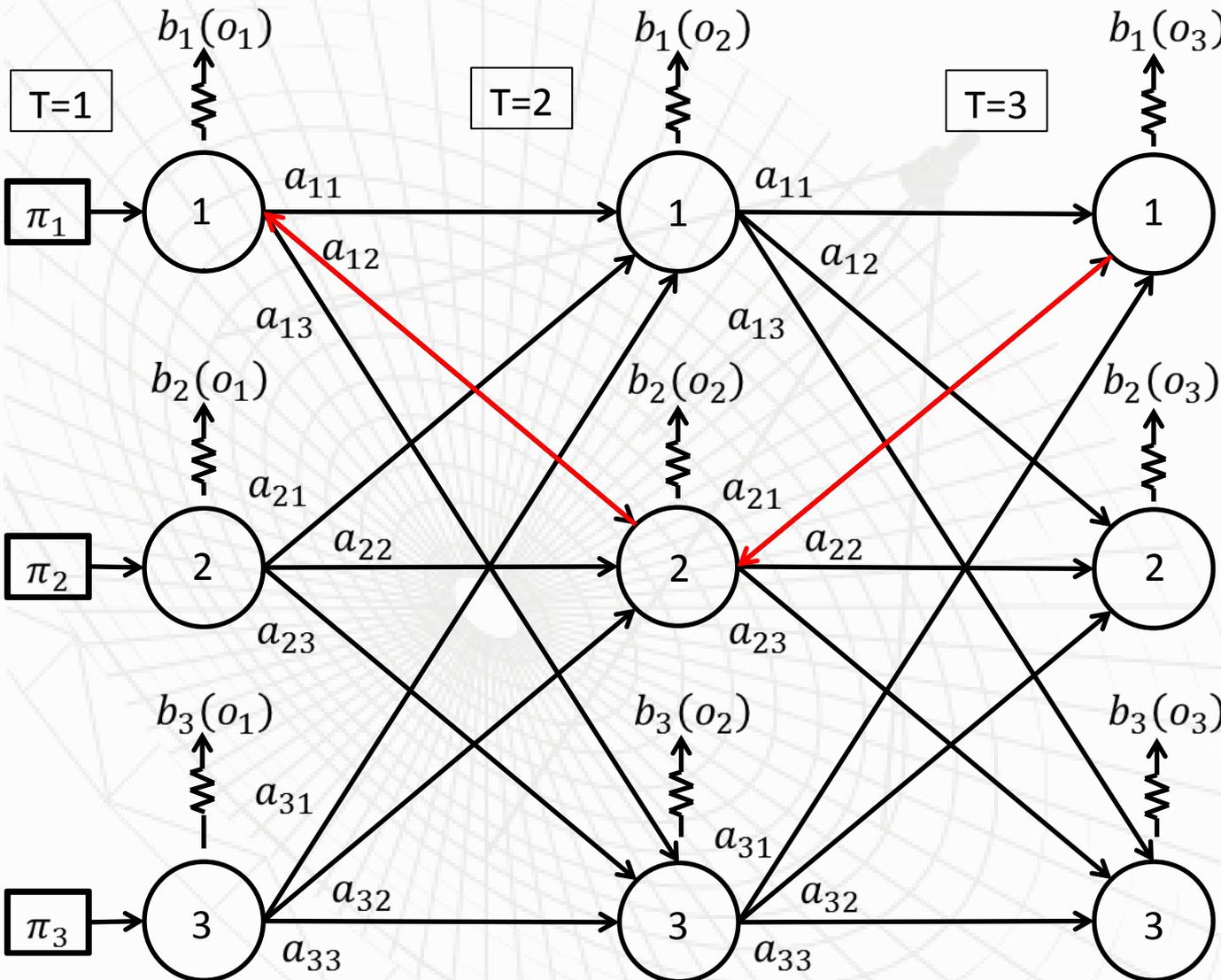
Apprentissage et optimisation architecturale de PROFHMM:



PROFHMM: RECONSTRUCTION



ALGORITHME DE VITERBI



Programmation Dynamique

Calcule:

- Probabilité maximale d'arriver à chaque état à chaque pas de temps.

Utilise:

- Calculs précédents.
- Probabilités du modèle.

Etape finale:

- Retrouve la série d'états cachés la plus probable.

- 
- *Introduction*
 - *La Méthodologie Statistique*
 - ***Applications***
 - *Complétion de données*
 - *Conclusions - Perspectives*

APPLICATION: BATS - DONNEES

PROFIL VERTICAUX:

NEMO – PISCES

Chlorophylle-A:
17 mesures verticales

Température:
9 mesures verticales

Période : 1992-2008
1241 pas de temps.

0-217 m



BATS (32°N – 64°W)

DONNÉES DE SURFACE:

SATELLITAIRES

Température de surface
Chlorophylle-A de surface

issues de **MODIS**
période 2002-2008:
402 pas de temps.
Données manquantes:
Interpolation par splines

NEMO

- Elévation du niveau de la mer
- *Radiation à faible longueur d'ondes (forçage)*
- *Intensité du vent (forçage)*

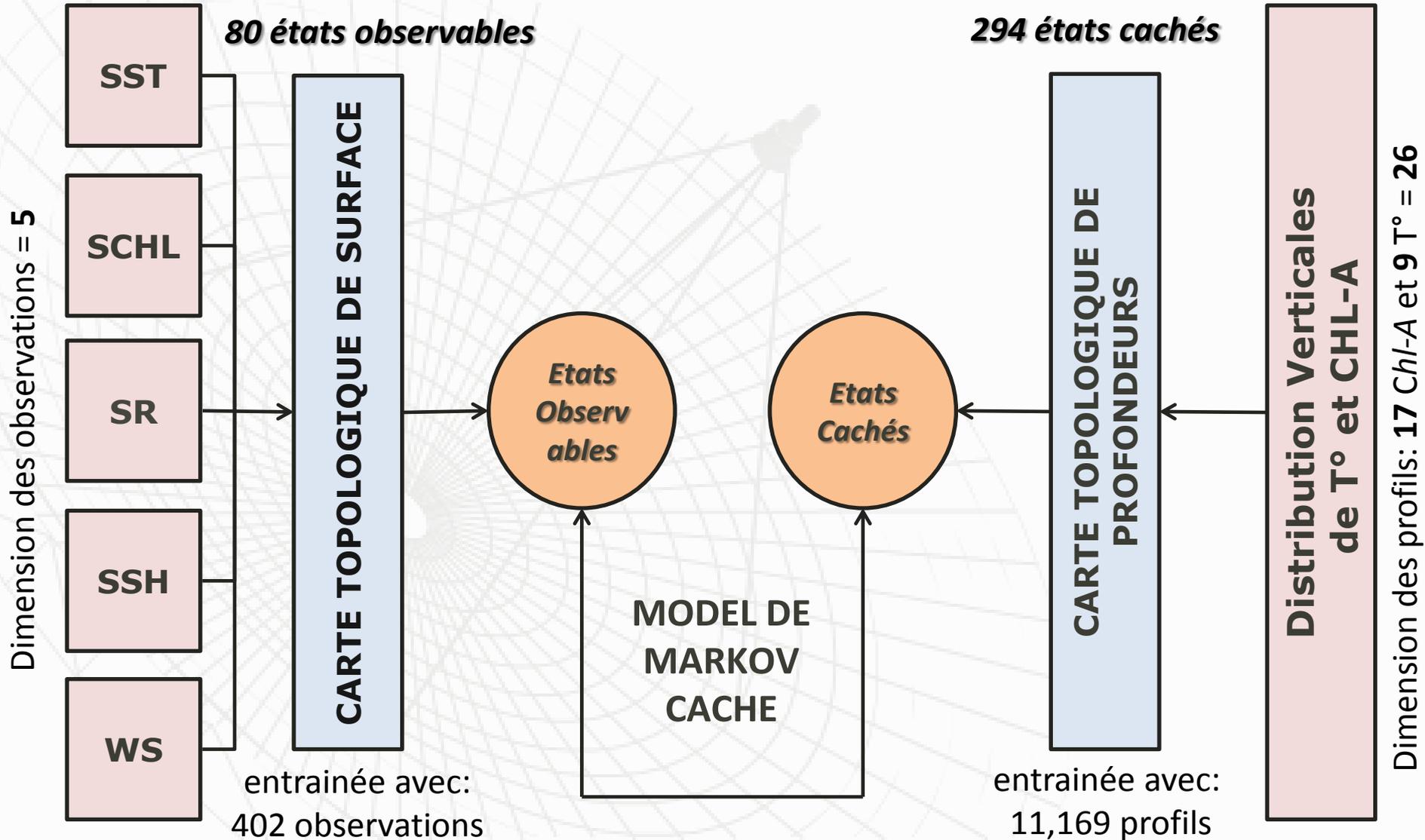
Moyennes sur 5 jours.

*Sélection des variables
utilisées par une étude
en ACP*

Egalement Appliquée:

*Période 1992-2008
Apprise totalement sur
NEMO-PISCES*

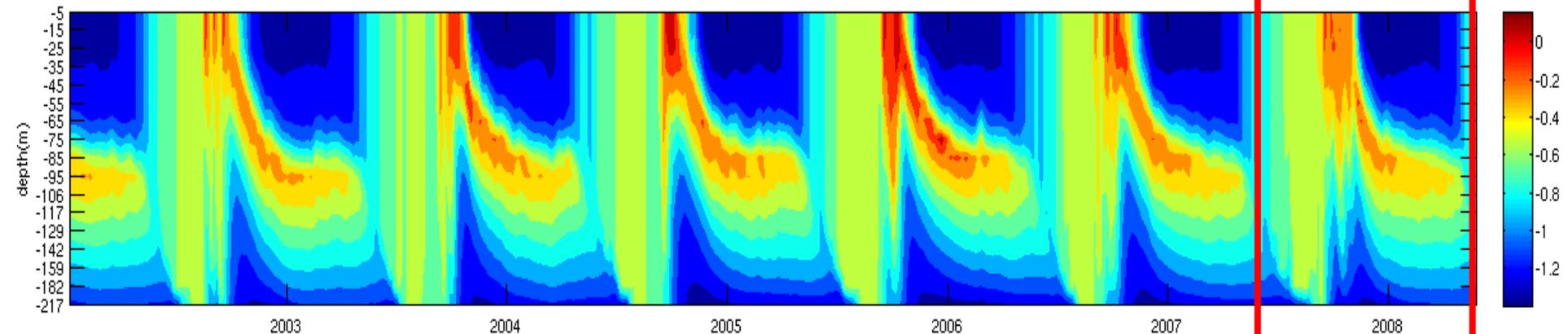
APPLICATION: BATS APPRENTISSAGE



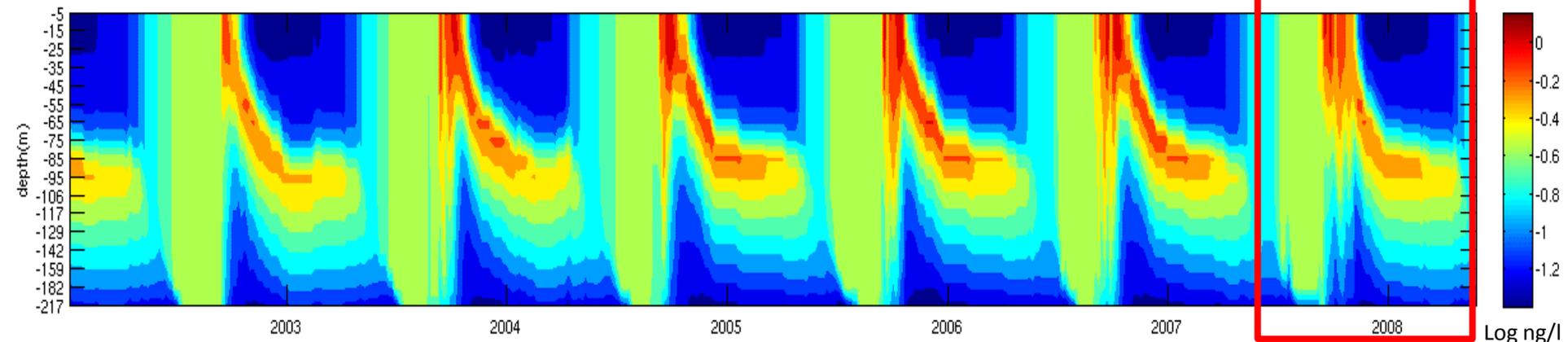
APPLICATION: BATS RECONSTRUCTION A partir d'Observations Satellite

Valeurs du model NEMO-PISCES

Validation



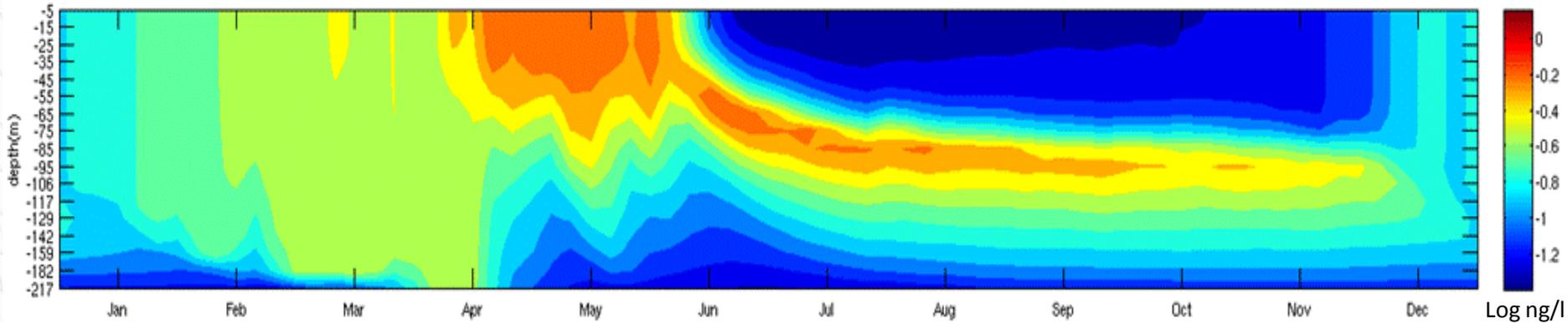
Reconstruction à partir de la surface avec PROFHMM



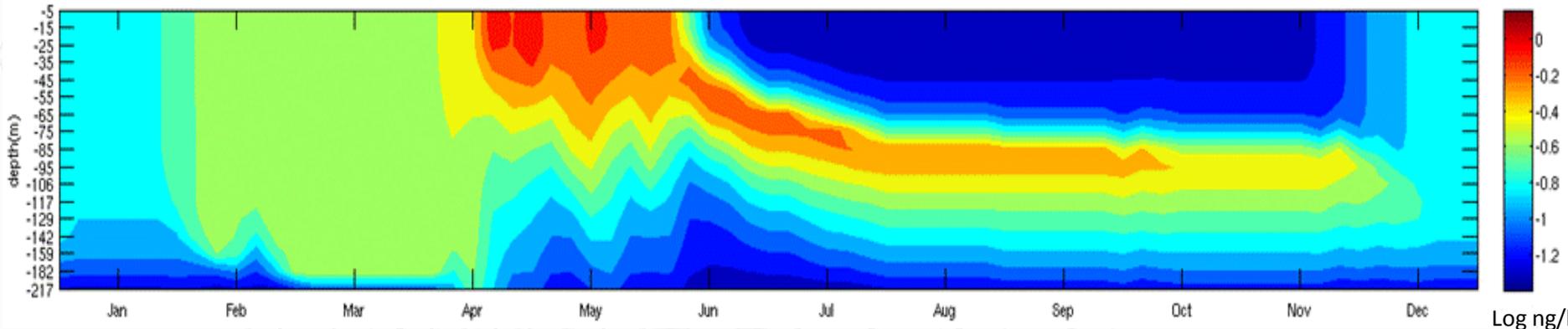
APPLICATION: Validation 2008

RMS Point / Point	10% plus faibles	10% plus fortes	AVERAGE
	0.0399	0.0096	0.0408

Valeurs du model NEMO-PISCES



Reconstruction à partir de la surface avec PROFHMM



APPLICATION: APOTRE

Analyse et Prévision Océanique de la Transparence des Eaux



APPLICATION: APOTRE

Données Observables

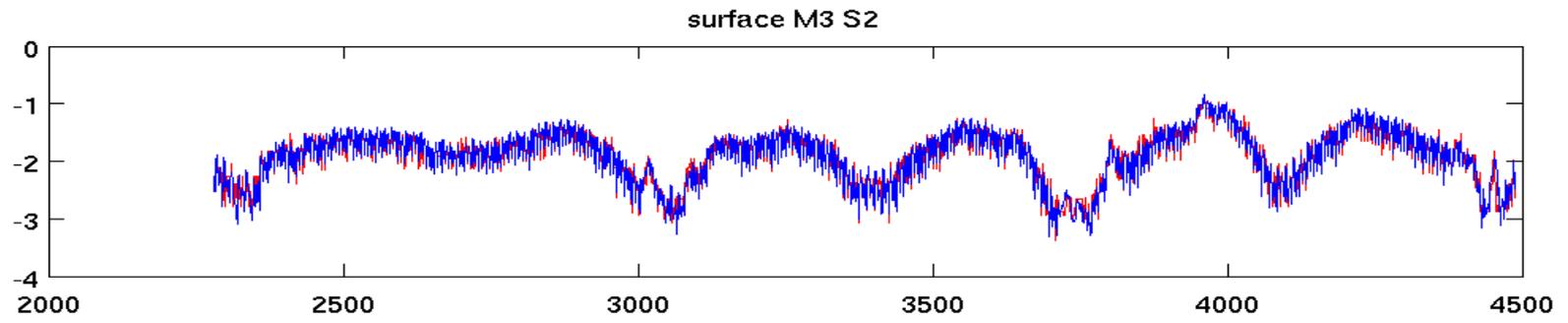
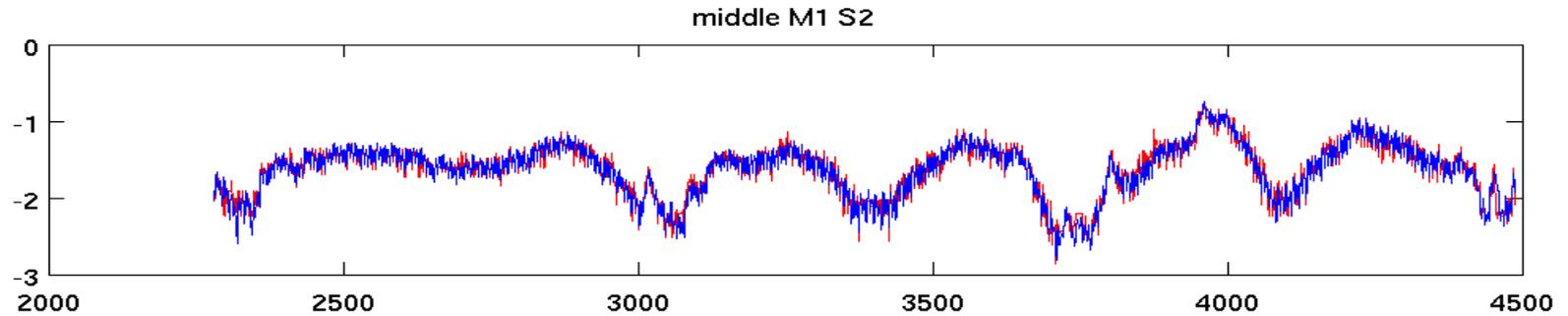
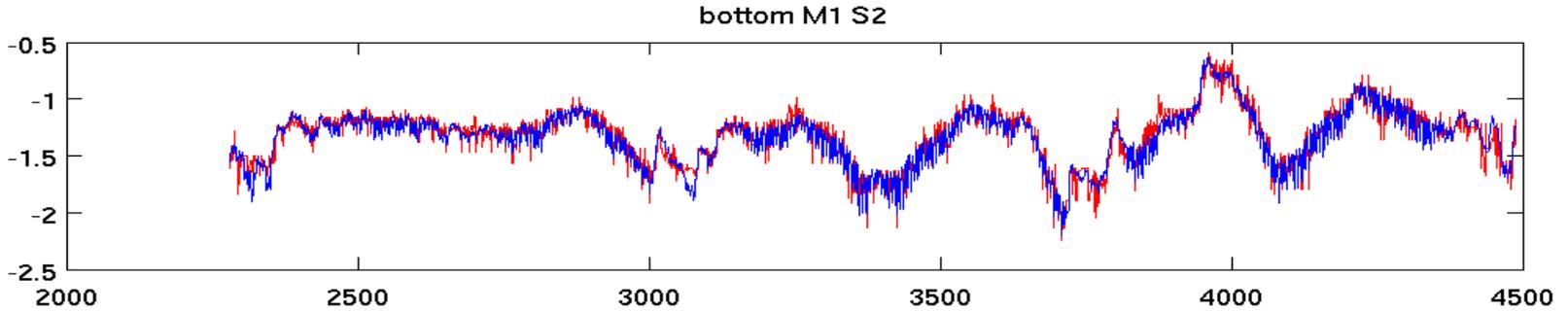
- Courants de marées (modèle numérique dédié): 2 directions (u,v), 11 niveaux de profondeur. (Réduction en 10 valeurs par ACP)
- Date et heure (en sinus, cosinus)
- Période des vagues.
- Intensité du vent.
- Tension au fond.
- Données Satellite:
 - Concentration totale de particules en surface
 - Gradient de concentration de particules en surface



Données Cachées

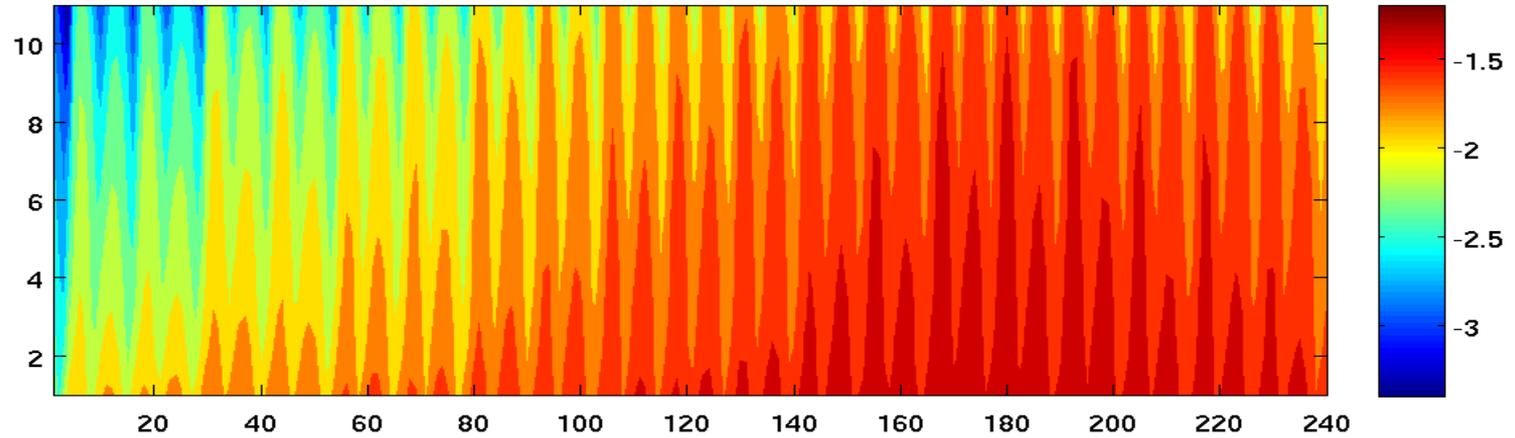
- 7 types de sables
- 11 niveaux de profondeur

APPLICATION: RECONSTRUCTION APOTRE

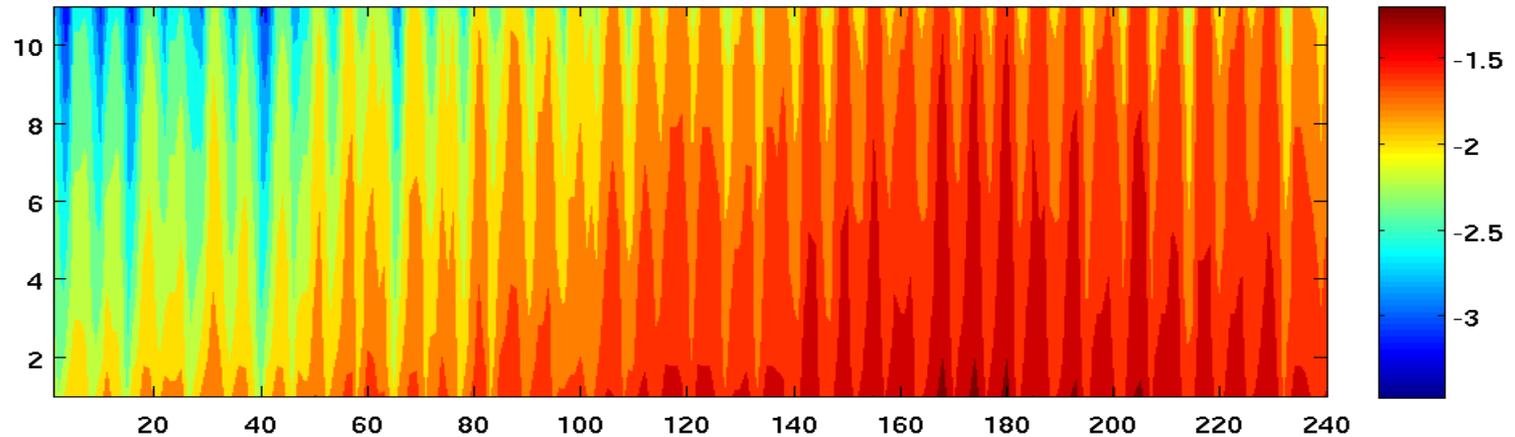


APPLICATION: RECONSTRUCTION APOTRE

ROMS days 516 525 M1 S2



PROFHMM days 516 525 M1 S2

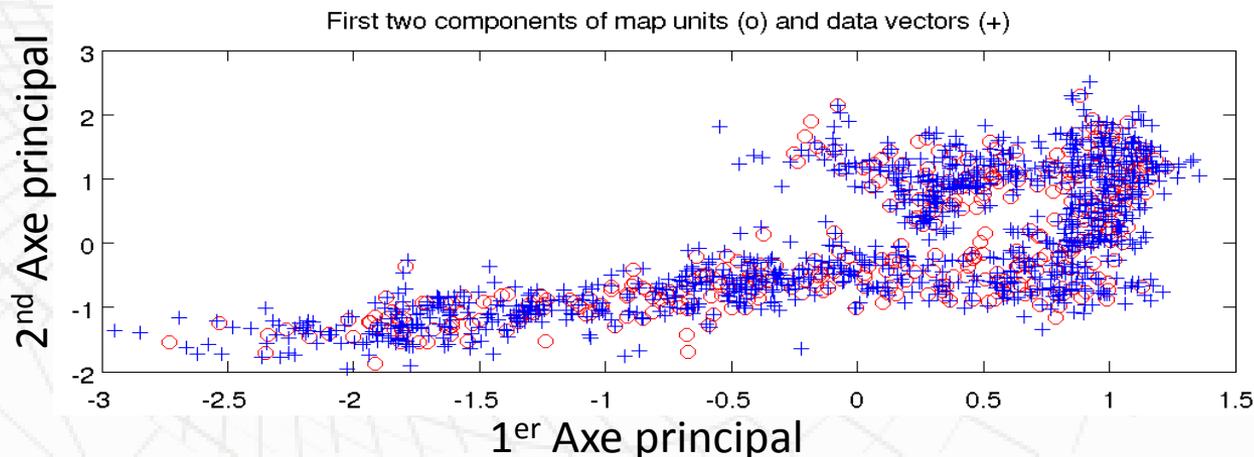


PROFHMM

- ***Permet de synchroniser la dynamique profonde avec la dynamique de surface.***
- Permet la génération d'un modèle statistique en synchronisant un modèle numérique avec des données satellite
- Peut générer un modèle statistique à partir de données in-situ.
- 1D temporel ou spatial: reconstruction de formes et intensités cohérentes.
- Assez robuste au bruit.
- Supporte grandes dimensionnalités.

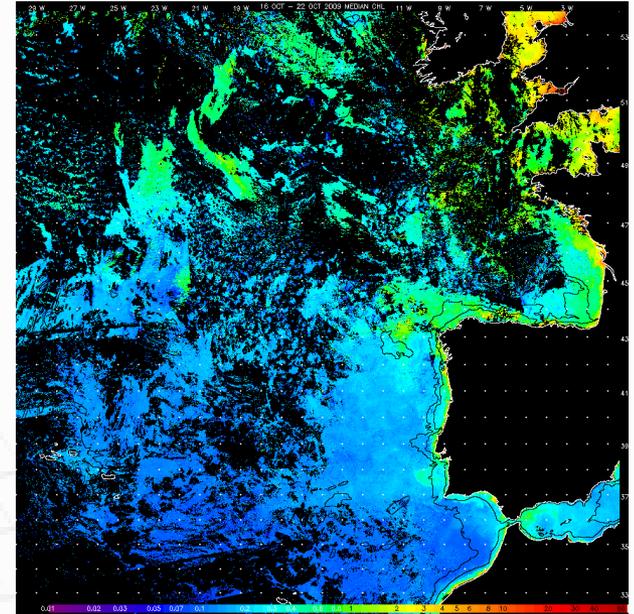
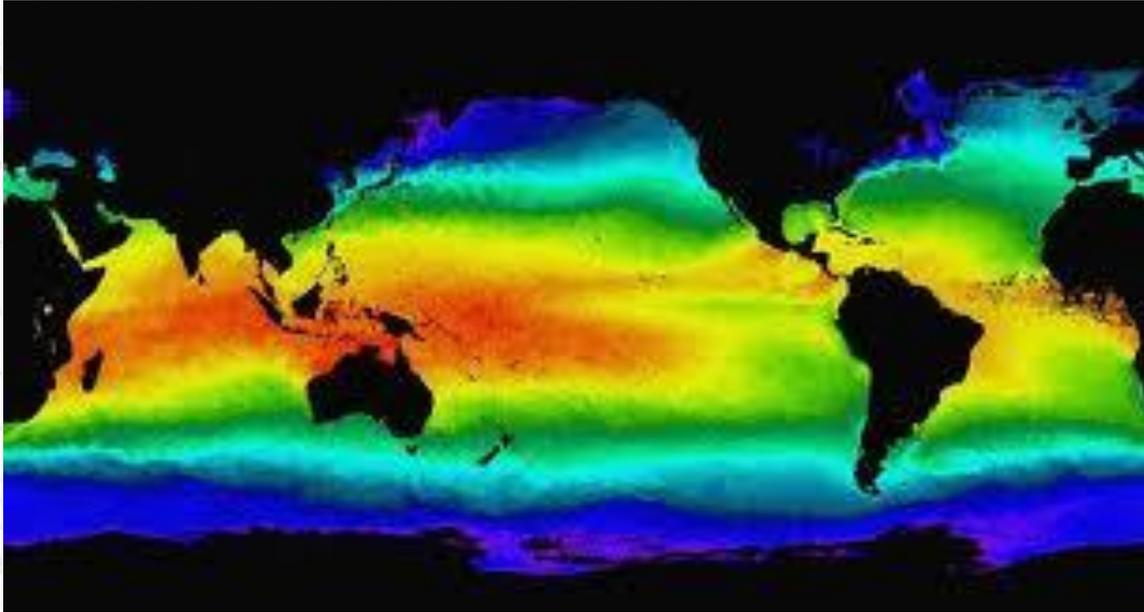
Mais:

- Impossible de reconstruire situations non observées/ pas couvertes par le modèle (mais référents extrêmes peuvent indiquer des situations extrêmes).



- 
- *Introduction*
 - *La Méthodologie Statistique*
 - *Applications*
 - ***Complétion de données***
 - *Conclusions - Perspectives*

DONNEES DE SURFACE INCOMPLETES



Données de surface utilisées:

- Sorties de Modèles Numériques

Ou si satellitaires:

- Produits de Réanalyses
- Incomplètes / Bruitées

PRECEDEMMENT:

***INTERPOLATION
PAR SPLINES***

PROFHMM_UNC

- Méthodologie qui prend en compte un expertise externe sur l'incertitude des données observées, et modifie l'algorithme de Viterbi pour moins pondérer le chemins vis-à-vis des observations.
- Application sur APOTRE, amélioration des résultats obtenus de l'ordre de $\sim 3\%$
- Dépend de la façon de sélectionner la certitude aux données.

ITCOMP SOM

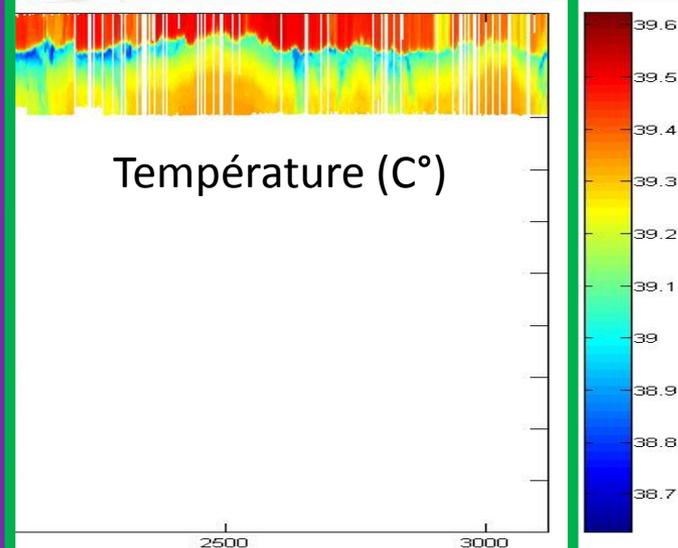
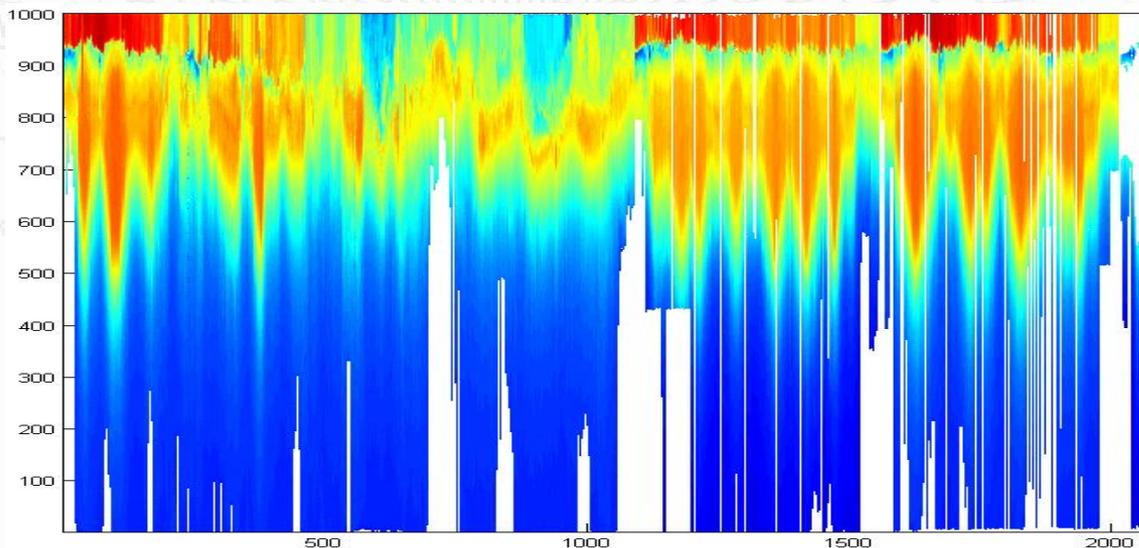
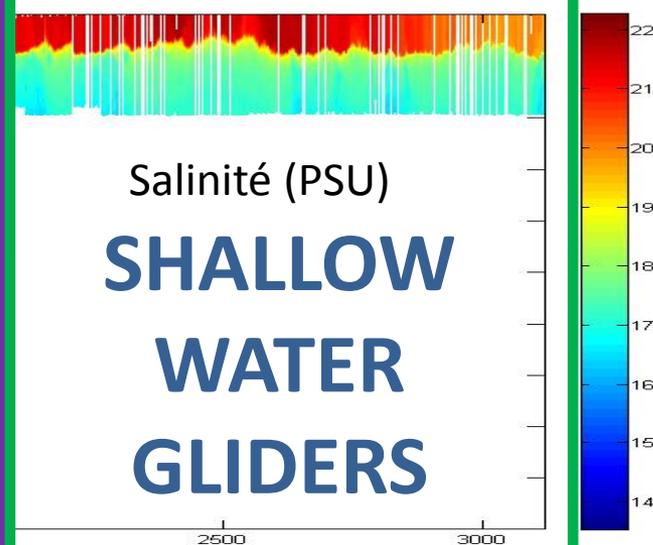
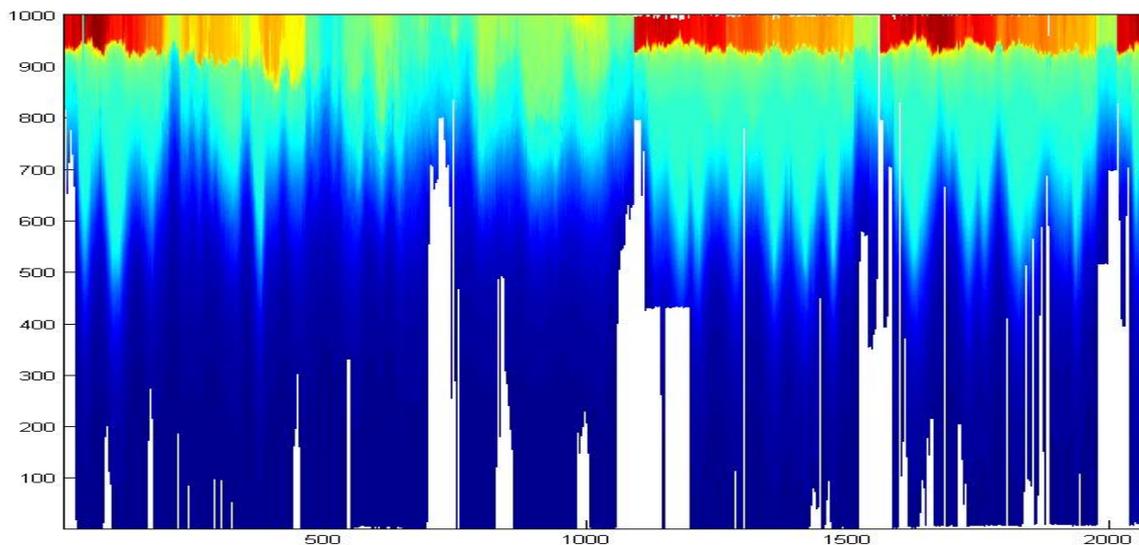
(Iterative Completion SOM)

- Complétion de bases de données contenant des données *très corrélées*, contenant un nombre important de *variables manquantes*.
- Utilise la corrélation locale des *variables présentes* avec les *variables manquantes* pour pondérer le calcul de distance entre les vecteurs et les référents.
- Complète et incorpore *itérativement* des données de plus en plus incomplètes.

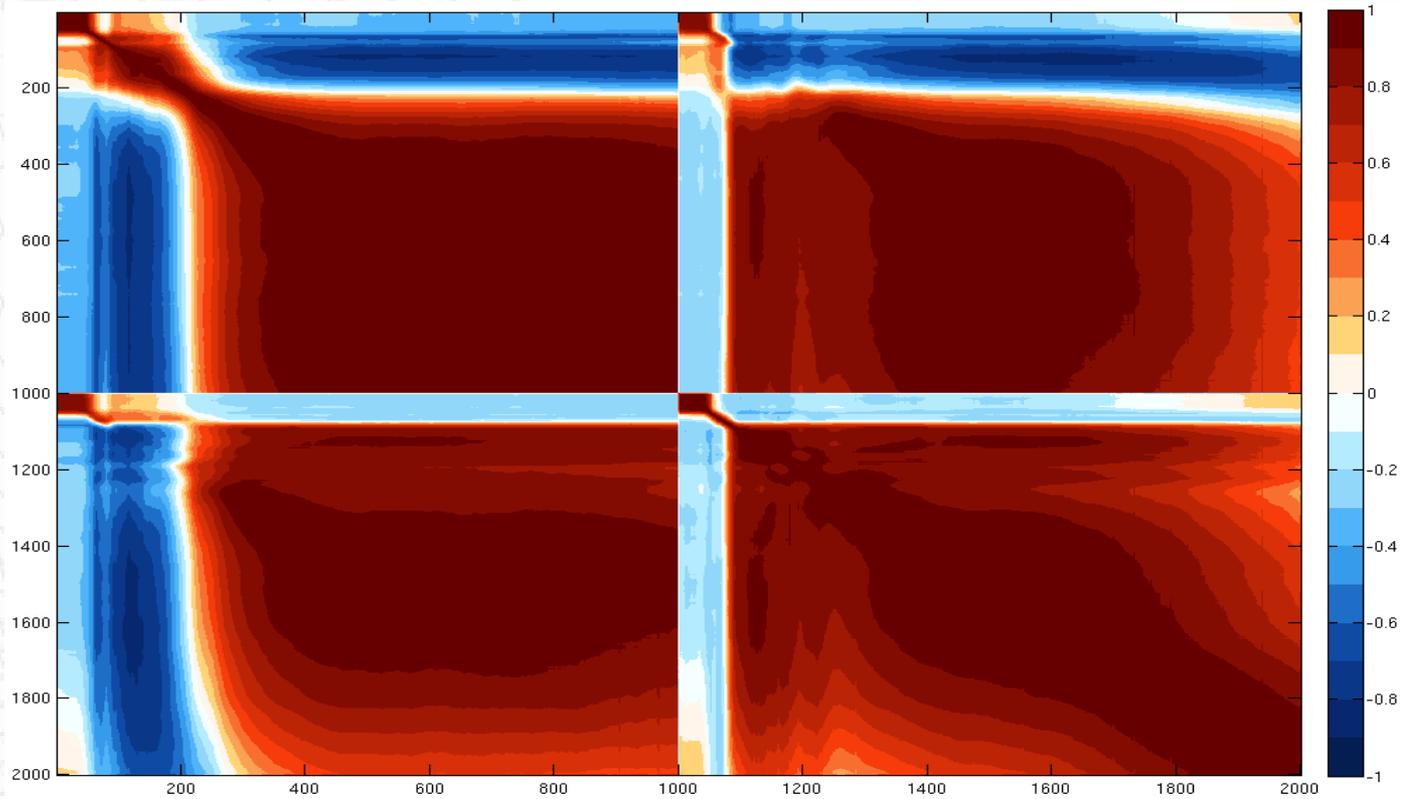
33% de données manquantes.

Eye of the Levantine

DEEP-SEA GLIDERS

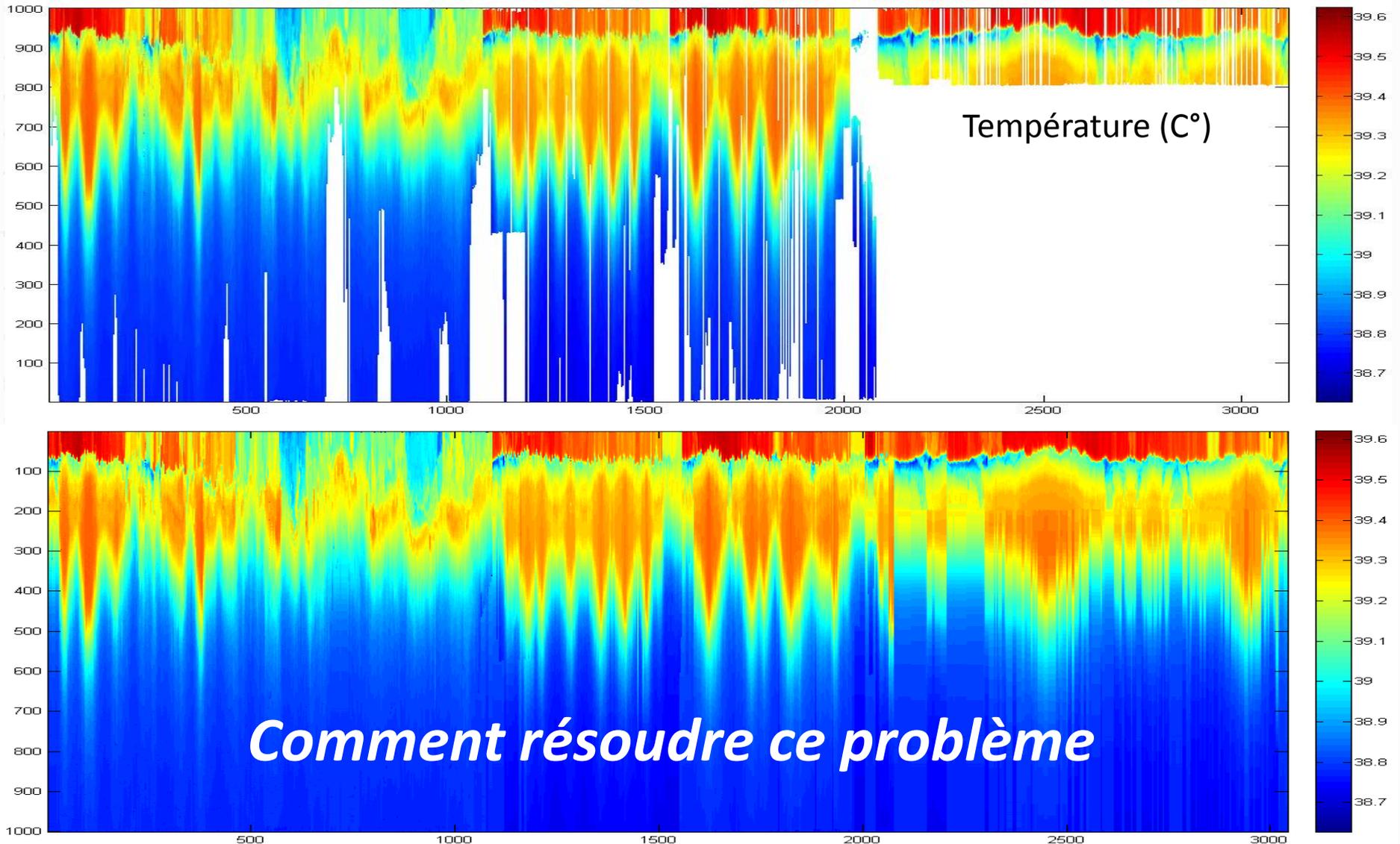


Eye of the Levantine



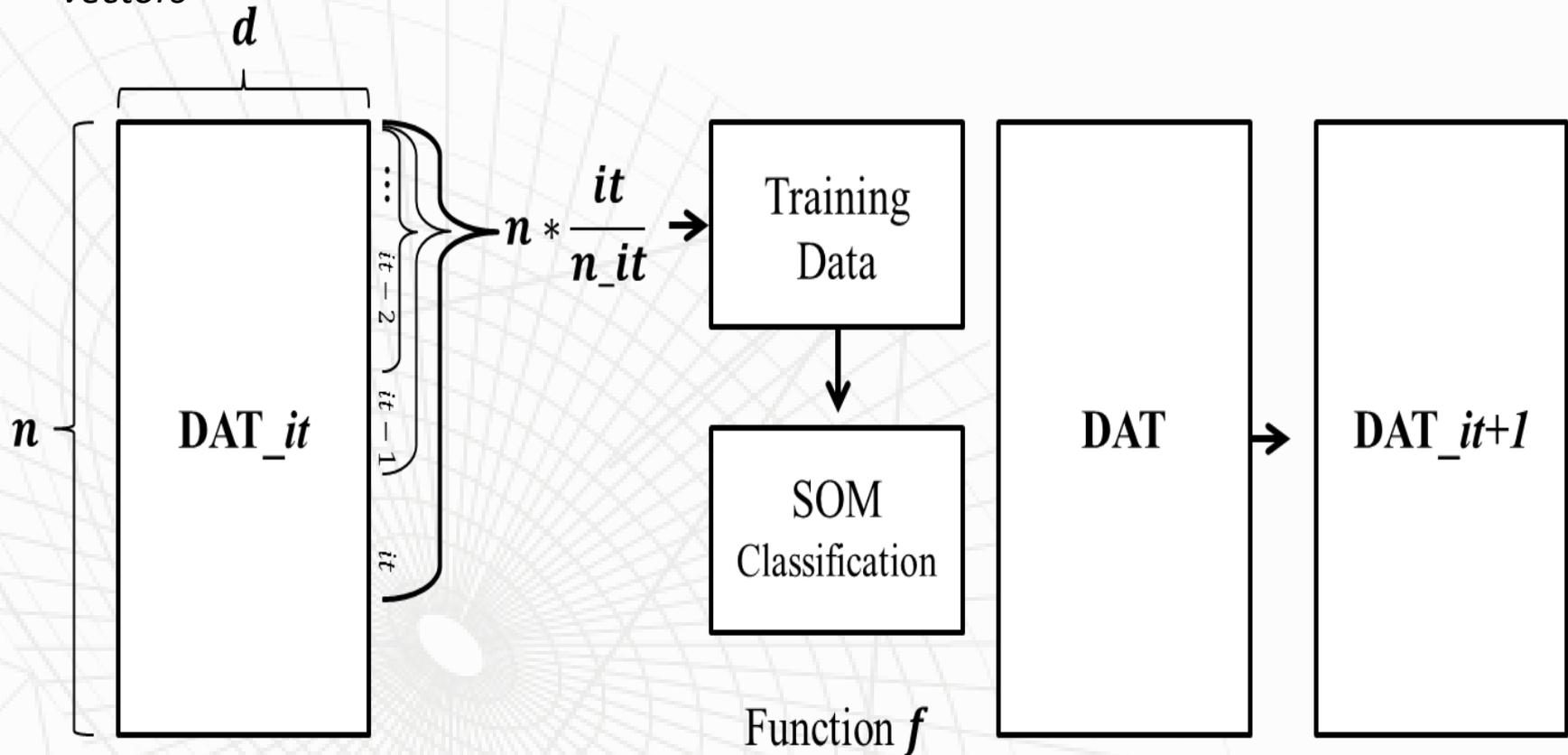
Correlogram des différents paramètres. Les 1000 premiers éléments correspondent à la température ($^{\circ}\text{C}$) dans la colonne d'eau allant jusqu'à 1000 mètres, les 1000 seconds correspondent à la salinité (PSU).

Eye of the Levantine



DAT is sorted by
"fullness" of
vectors

Entraînement Itératif



Initialisation:

$$DAT_1 = DAT$$

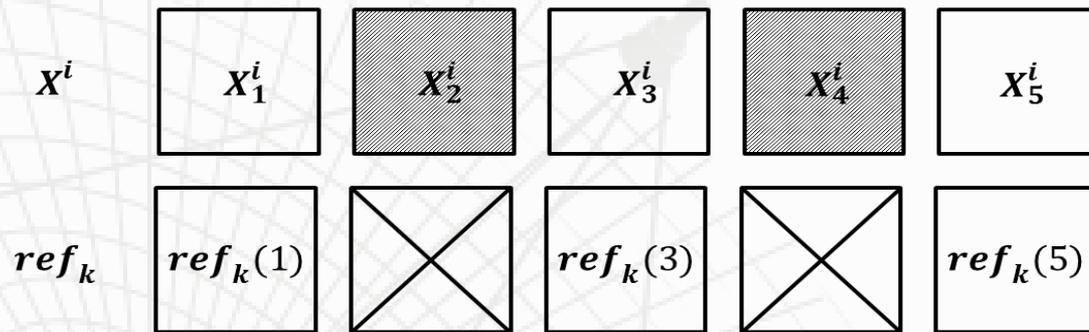
Used for the
completion of:

Selection du “Best-Matching Unit”

$$full_i = \{1, 3, 5\}$$

$$missing_i = \{2, 4\}$$

$$w_k(j) \propto \begin{cases} abs(covariance(j, 2)) \\ abs(covariance(j, 4)) \end{cases} \text{ around } ref_k$$

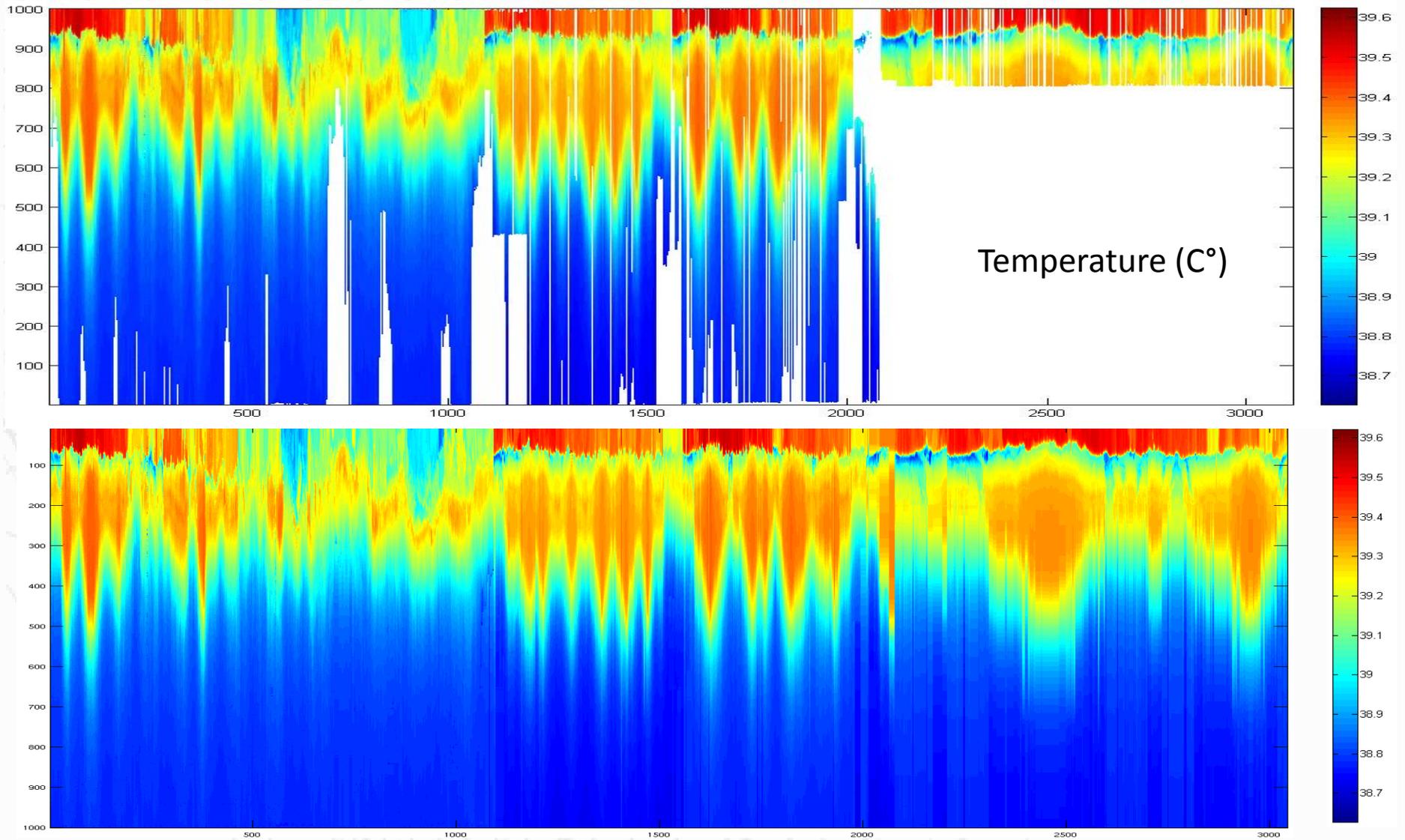


$$sim(ref_k, X^i) = \sqrt{\sum_{j \in \{1,3,5\}} ((ref_k(j) - X_j^i) * w_k(j))^2}$$

$$w_k(j) = 1 + \sum_{m \in missing_i} thr_k(m, j)$$

$$thr_k(m, j) = \begin{cases} 0, & \text{if } abscov_k(m, j) < tr_value \\ abscov_k(m, j), & \text{else} \end{cases}$$

Completion



ITCOMP SOM

- Aussi utilisé pour la complétion d'une base de données pour retrouver la $p\text{CO}_2$ en mer Baltique.
- 1445 vecteurs contenant huit paramètres (60% avec au moins une valeur manquante):
(SST SCHL CDOM NPP MLD DAY_{\sin} DAY_{\cos} $p\text{CO}_2$)

Nombre de Paramètres Manquants	R	$p\text{CO}_2$ RMSE $\mu = 351.06$, $\sigma = 112.11$
0	0.96	25.7 μatm
1	0.93	39.5 μatm
2	0.86	31.9 μatm
3	0.81	51.4 μatm

- Introduction
- *La Méthodologie Statistique*
- *Applications*
- *Complétion de données / Prise en compte d'incertitudes*
- **Conclusions - Perspectives**

CONCLUSIONS

- ***Outil d'inversion.***
- ***Reconstruction de formes et intensités cohérentes.***
- ***Prise en compte de grandes dimensions.***
- ***Prise en compte des données manquantes ou incomplètes avec PROFHMM_UNC / ITCOMPSOM.***
- ***Utilisations diverses:***
 - ***Simplifier un modèle numérique.***
 - ***Générer un modèle statistique à partir de données in-situ.***
 - ***Etudes de cas.***
- ***Facile à implémenter.***

PERSPECTIVES

Sur le plan *algorithmique*:

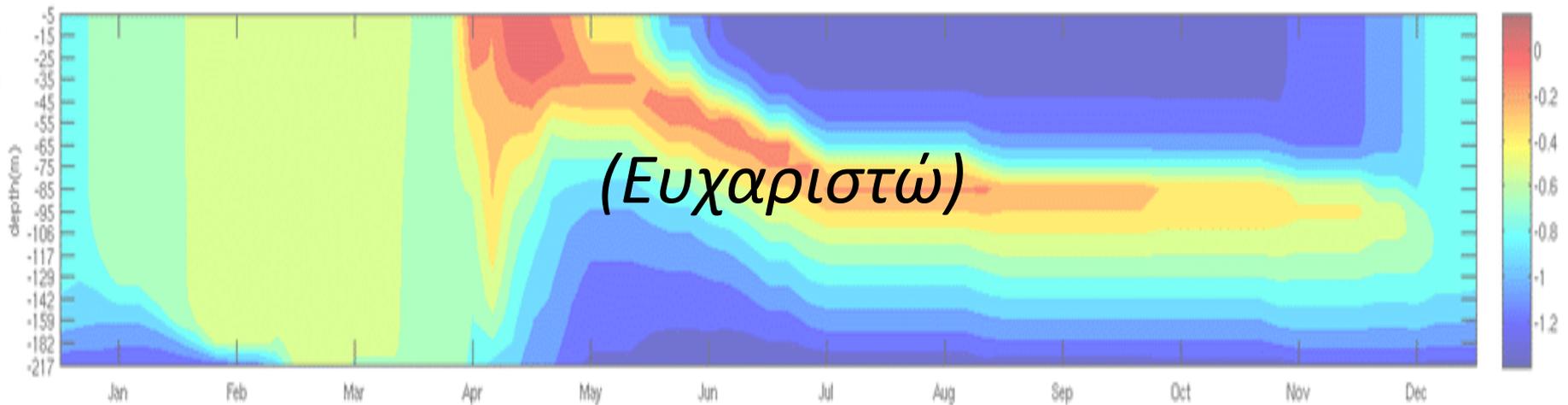
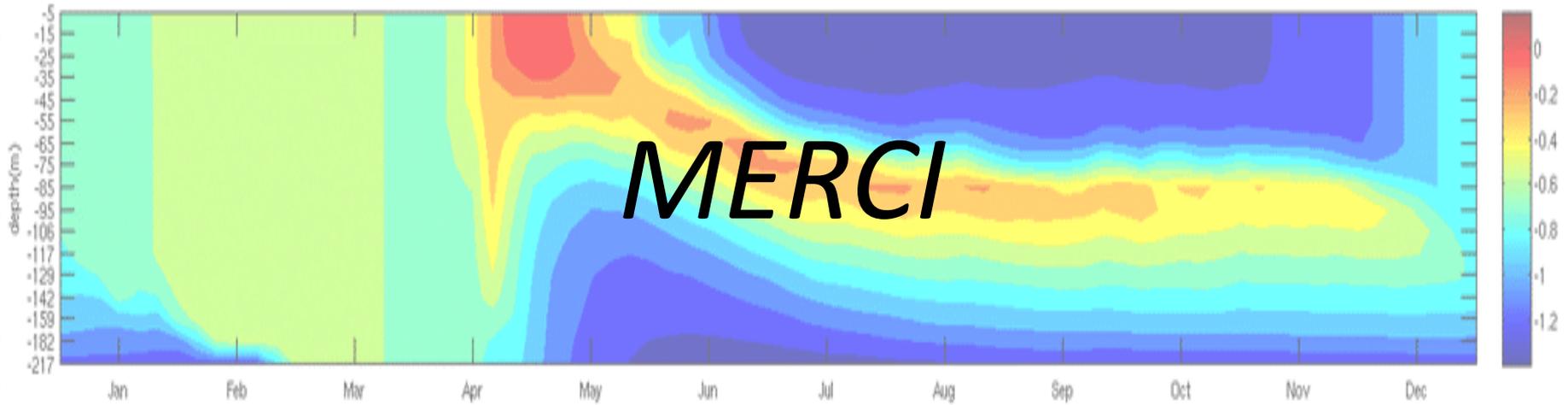
- Passage aux champs de Markov 2D d'observations (théorie faite, application en cours).
- Générateur de séries d'observations réalistes.
- Méthode de complétion de données.

Sur le plan *méthodologique* (contexte géophysique):

- Aide à l'assimilation:

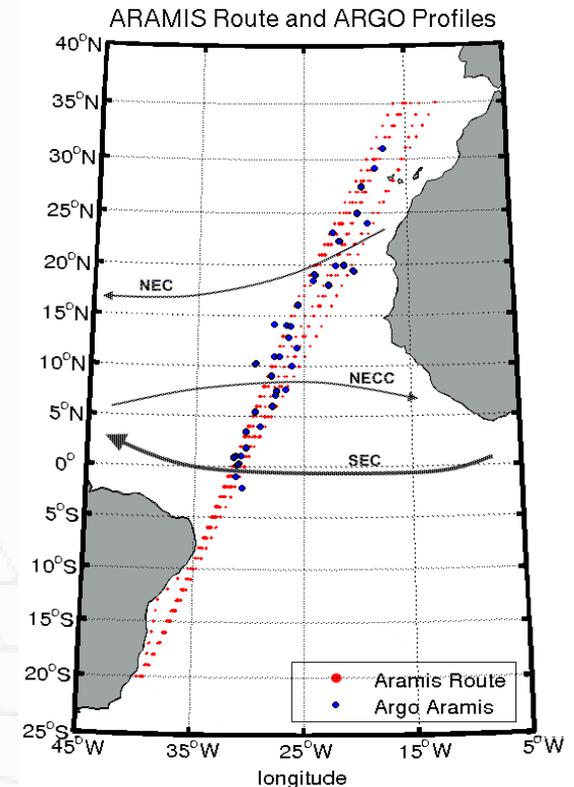
Détermination du first guess.

Année 2005 sur BATS



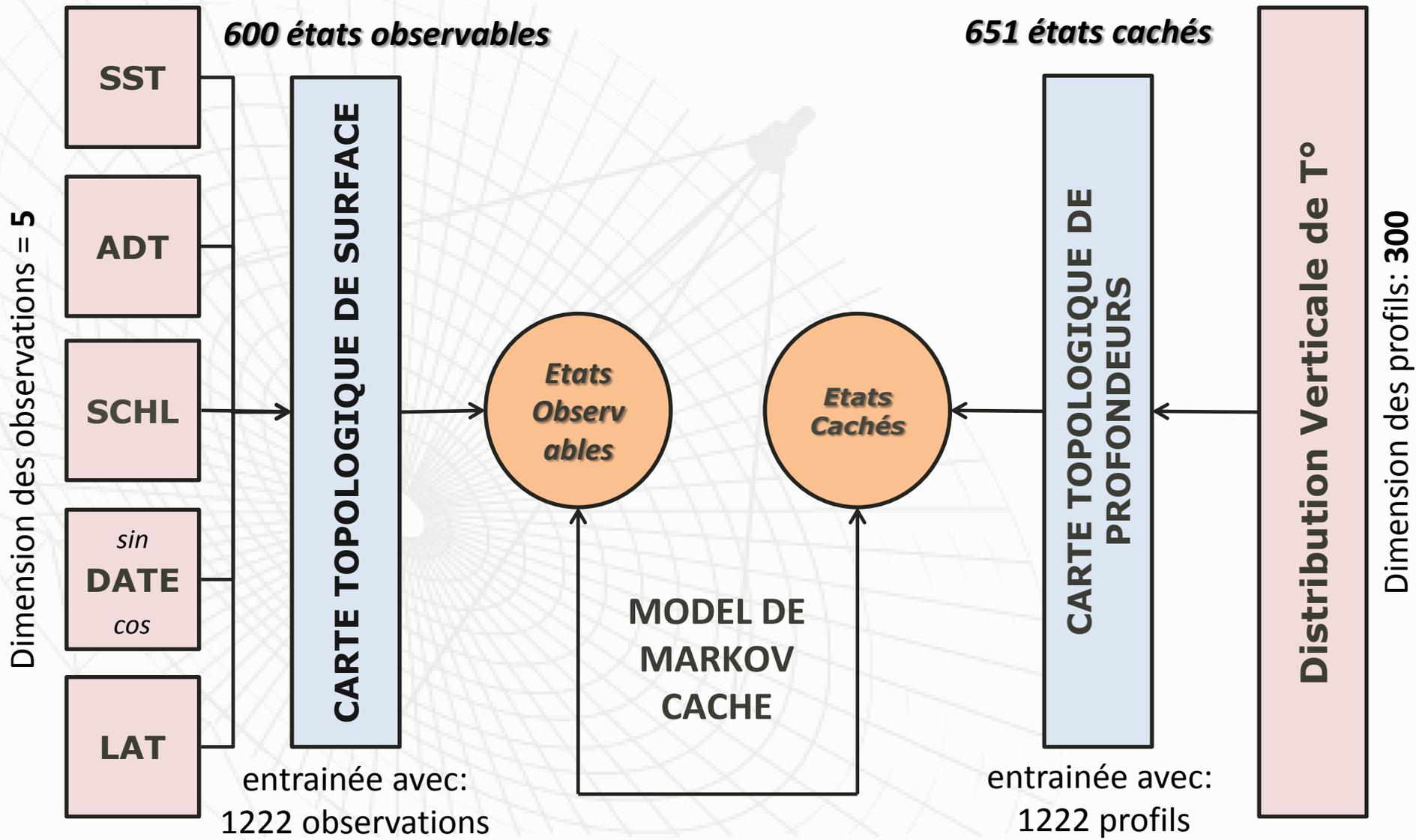
APPLICATION: ARAMIS

- ARAMIS (Altimétrie sur un Rail Atlantique et Mesures In Situ)
PI S.ARNAULT: Surveillance des structures thermo halines des couches de surface océaniques en Atlantique tropical entre 2002 et 2008.
- 1300 Mesures in-situ par XBT de T °C (utilisation de 5 à 305 mètres)
- Précision des mesures: +/- 0.1°C
- Mesures satellites de T° de surface, de la topographie dynamique absolue (ADT)



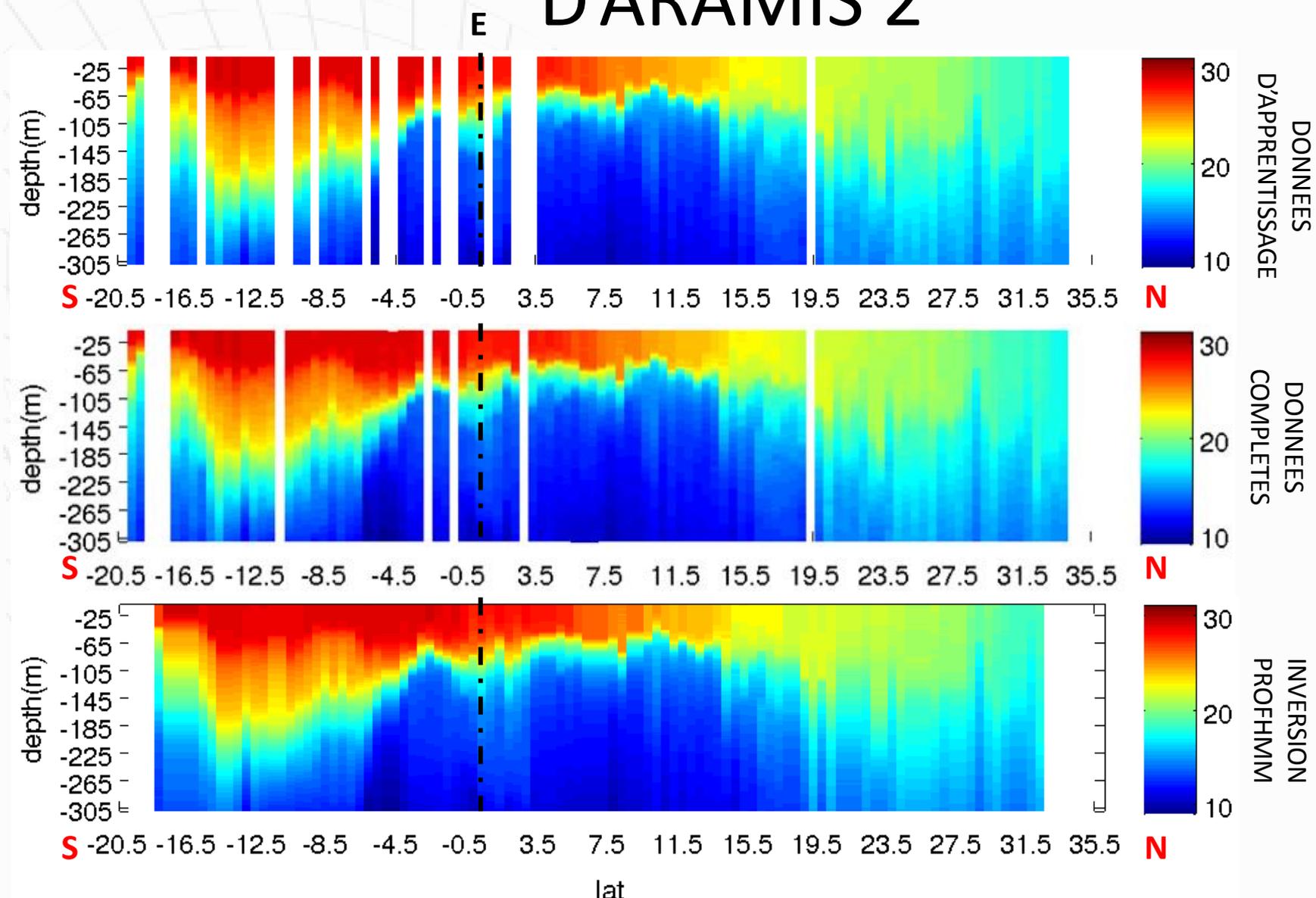
Hypothèse: Temps \Leftrightarrow Espace

APPLICATION: ARAMIS APPRENTISSAGE



APPLICATION: RECONSTRUCTION

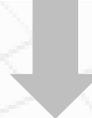
D'ARAMIS 2



(2003/03/15 – 2003/03/23)

RECONSTRUIRE EN UTILISANT LA DYNAMIQUE DE SURFACE

**Contraintes de la dynamique de
surface**



**Complétion
de séries de
données de
surface**

Contraintes des observations



**INDICATEUR DE CONFIANCE DES
OBSERVATIONS**

MODIFICATION DE L'ALGORITHME DE VITERBI

Alternative:

PRISE EN COMPTE DES INCERTITUDE:

***MODIFICATION DE
L'ALGORITHME DE VITERBI***

$conf(Obs) \in [0,100]$

Fonction de confiance



Expertise de qualité
de l'observation.

$F_w(X, Y) \in \{[0,1], [0,100]\}$

Fonction de pondération



Modifié les probabilités
d'émission en fonction de la
confiance

EXPERIENCE JUMELLE

NEMO-PISCES SURFACE

(SCHL, SST, SSH, SR, WS)



Etats cachés du Modèle de Markov Caché

**SIMULATION DE
COUVERTURE
NUAGEUSE:**

sous - échantillonnage



*Etats observables du
Modèle de Markov Caché*

**FONCTION DE CONFIANCE
DEPENDANTE DU
SOUS-ECHANTILLONNAGE**

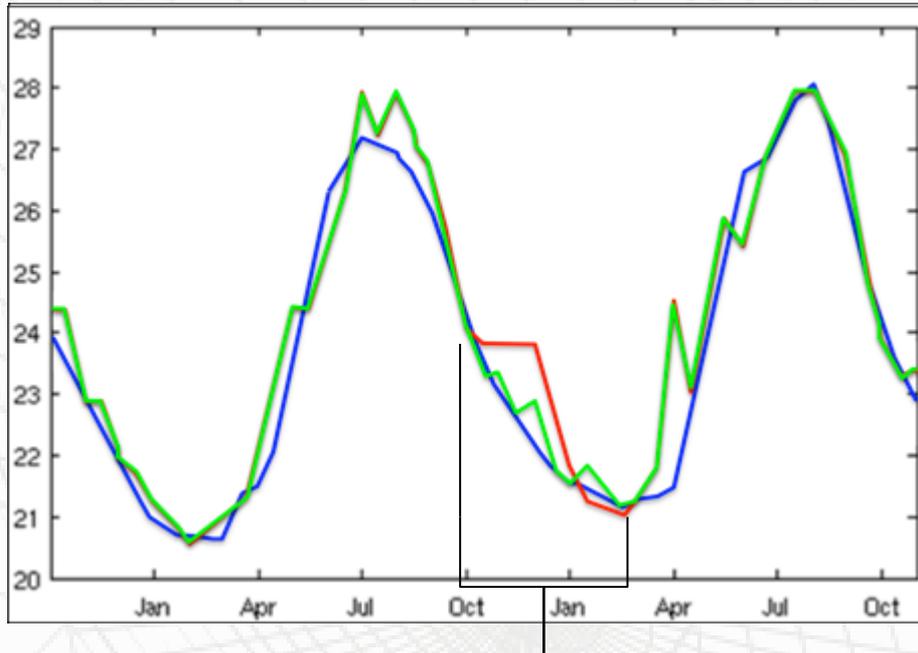


NEMO-PISCES SURFACE + BRUIT GAUSSIEN

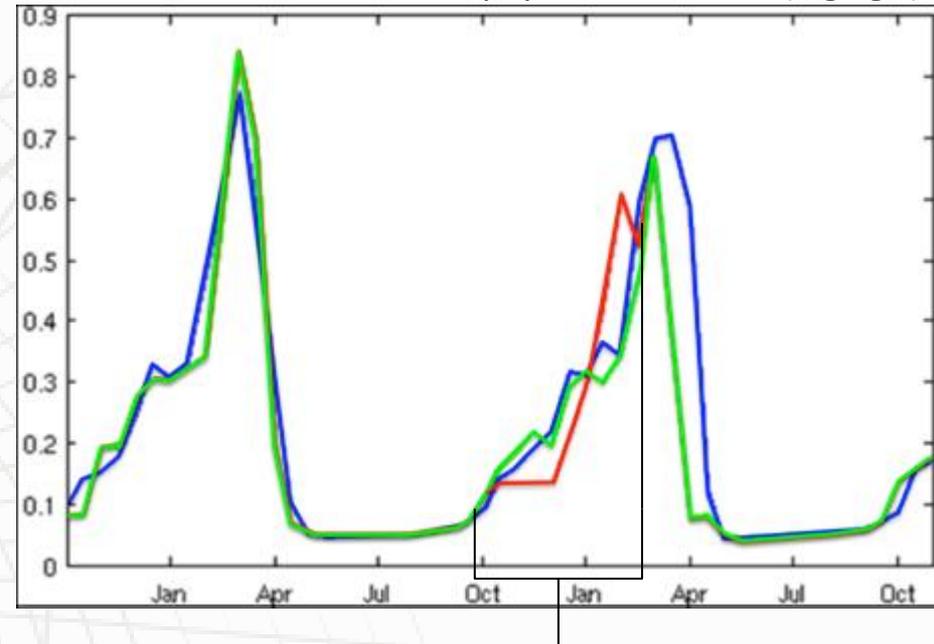
ETUDE DE SENSIBILITEE

Reconstruction de la SST SCHL à BATS 2007 -2008

Température de surface (°C)



Chlorophylle-A de surface (log ng/l)



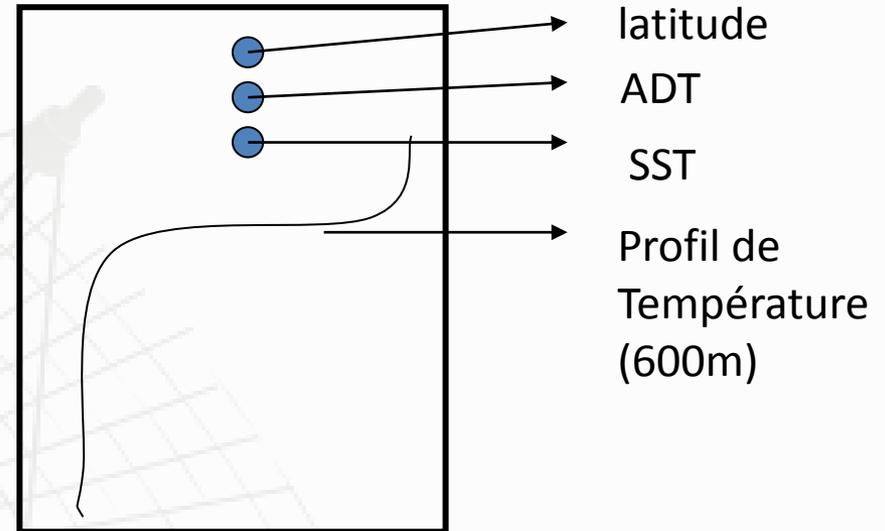
Période aux données de surfaces modifiées

- Données de surface NEMO (« REEL »)
- Reconstruction PROFHMM de la surface à partir de obs sans tenir compte de la confiance
- Reconstruction PROFHMM_UNC en tenant compte de la confiance

Yves Tanguys

- Principe général

- Inversion d'une section entière à un t donné
- Respecter la cohérence spatiale de T



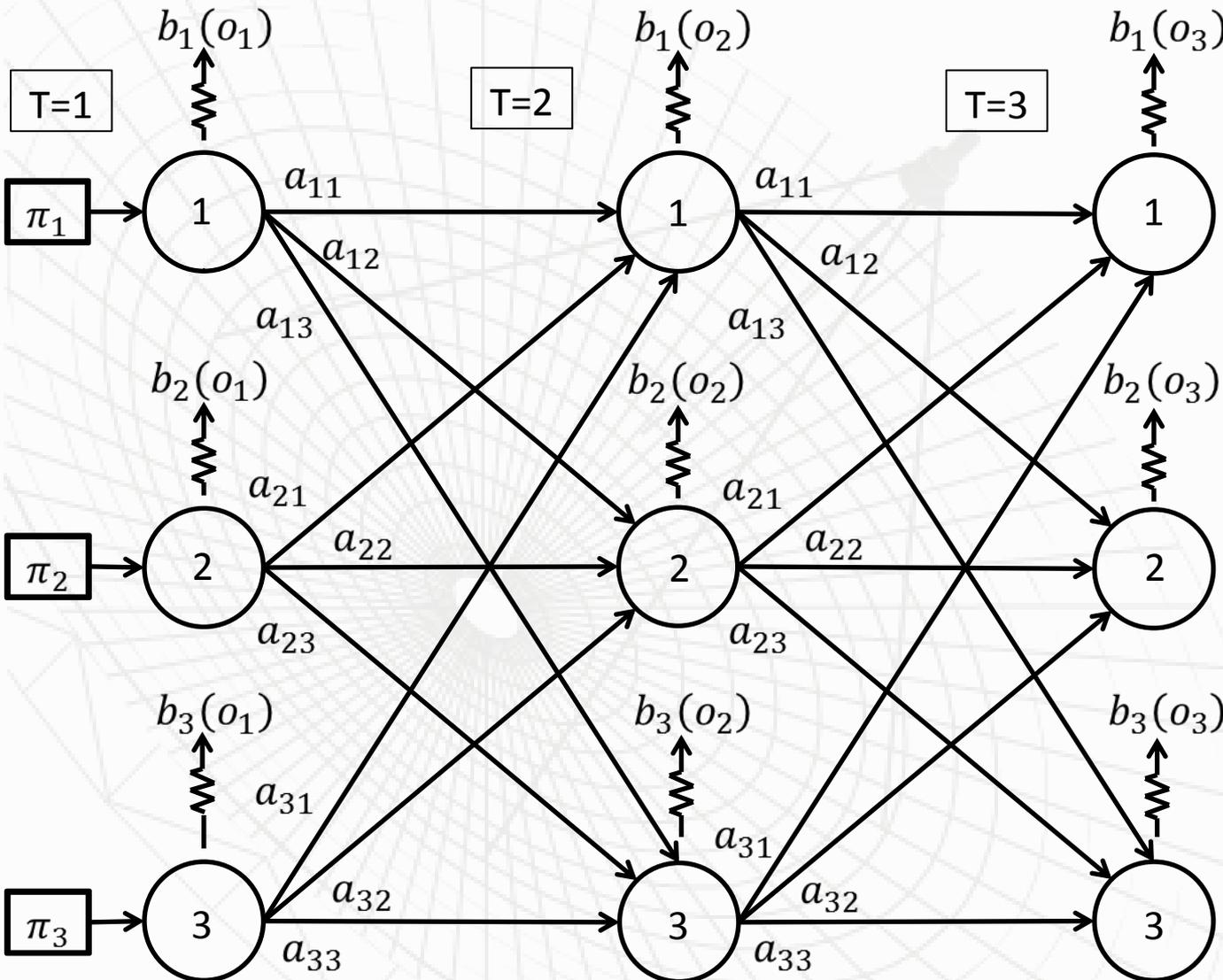
- Partie 1

- Obtenir un certain nombre de classe cohérentes en terme de donnée de surface

- Partie 2

- Choisir parmi elles la meilleur a priori en fonction des points proches

ALGORITHME DE VITERBI (1/3)



EVALUATION DIRECTE:

$(2T-1) * N^T$
Multiplications

$N^T - 1$
Additions

Soit:
 $\approx 2T * N^T$
opérations

PROFHMM_UNC: FONCTIONS DE CONFIANCE ET DE PONDERATION

The Viterbi Algorithm

Initialization:

$$\text{For } 1 \leq i \leq N_{hid}$$

$$\delta_1(i) = \pi_i * b_i(o_1),$$

$$\psi_1(i) = 0$$

with π_i the initial probabilities the state i .

Iterative calculation

For $2 \leq t \leq T$

For $1 \leq j \leq N_{hid}$

$$\delta_t(j) = \left(\max_{1 \leq i \leq N_{hid}} [\delta_{t-1}(i) * a_{ij}] \right) * b_j(o_t)$$

$$\psi_t(j) = \text{arg max}_{1 \leq i \leq N_{hid}} [\delta_{t-1}(i) * a_{ij}]$$

Ending values

$$P = \max_{1 \leq i \leq N_{hid}} [\delta_T(i)]$$

$$q_T = \text{arg max}_{1 \leq i \leq N_{hid}} [\delta_T(i)]$$

Backpropagating

For $t=T$ to 2

$$q_{t-1} = \psi_t(q_t)$$

Introduction de deux fonctions

$$conf(Obs) \in [0,100]$$

Fonction de confiance



Expertise de
qualité de l'obs.

100% confiance dans l'Obs



$$conf(Obs) = 100$$

0% confiance dans l'Obs



$$conf(Obs) = 0$$

$$\delta_t(j) = \left(\max_{1 \leq i \leq N} [\delta_t(i) * a_{ij}] \right) * F_w(b_i(o_t), conf(o_t))$$

$$F_w(X, Y) \in \{[0,1], [0,100]\}$$

Fonction de pondération

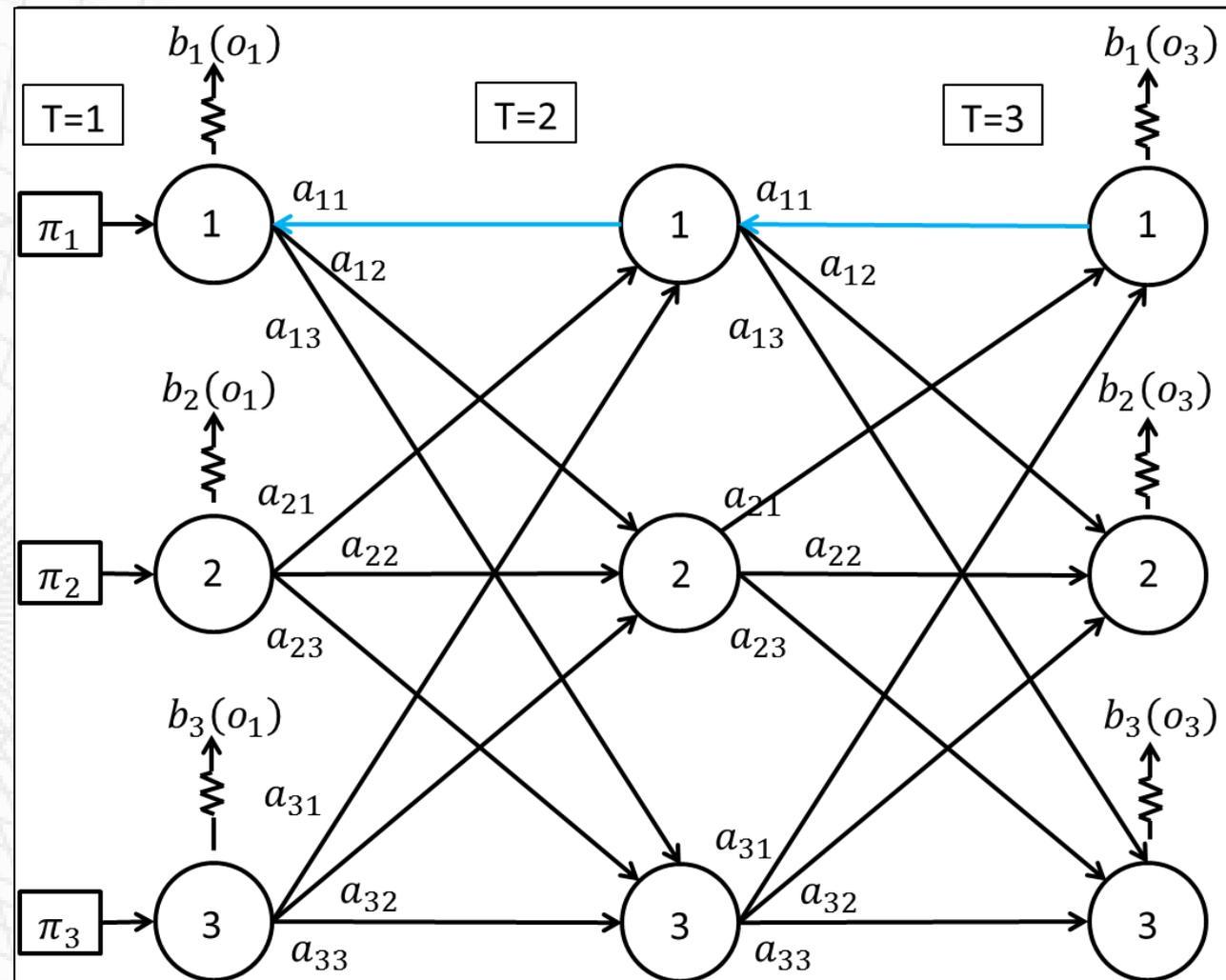
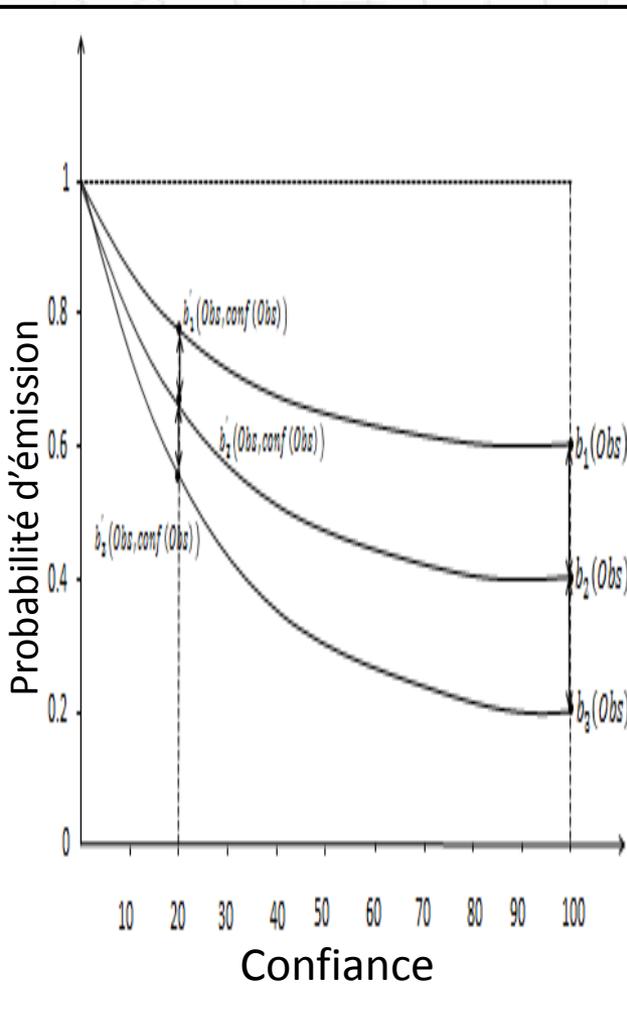


Monotone
envers X et Y
Décroissante

$$F_w(b_i(Obs), 100) = b_i(Obs)$$

$$F_w(b_i(Obs), 0) = 1$$

PROFHMM_UNC: CALCUL DU CHEMIN



ETUDE DE SENSIBILITE: DONNEES

SST

SCHL

SR

SSH

WS

Pour tester la Modification de VITERBI:

ETATS CACHES:

Données de surface de NEMO-PISCES à BATS

MOYENNE DES 3 PAS DE 5 JOURS

ETATS OBSERVES:

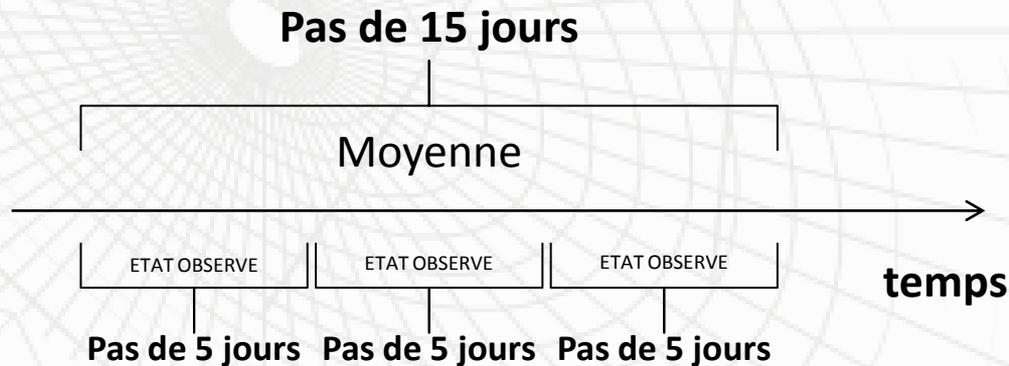
Données de surface de NEMO-PISCES à BATS

+ Introduction d'un bruit gaussien de variance $0.35 * \sigma$

MOYENNE DES 3 PAS DE 5 JOURS



BATS (32°N -64°W)

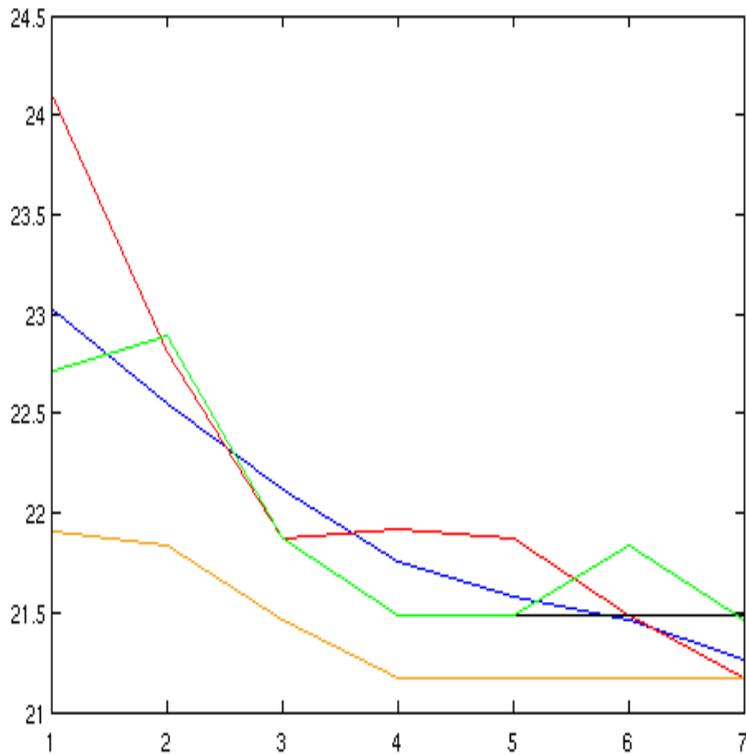


Observations Bruitées:
Que une données sur trois

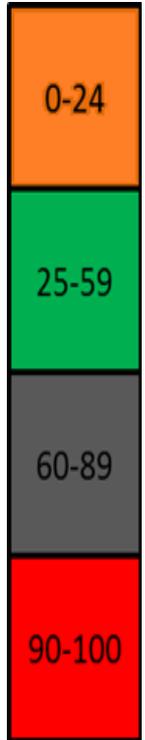
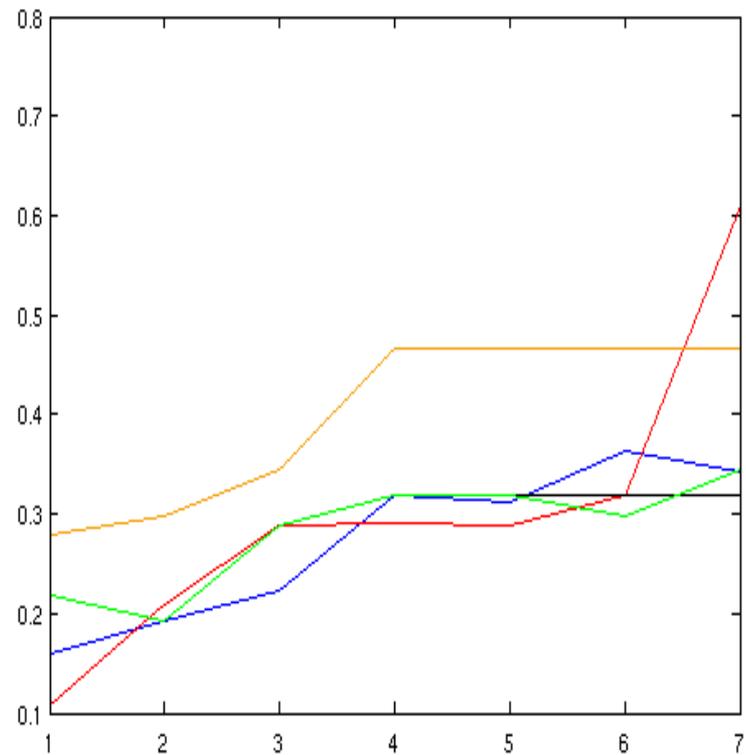


ETUDE DE SENSIBILITEE

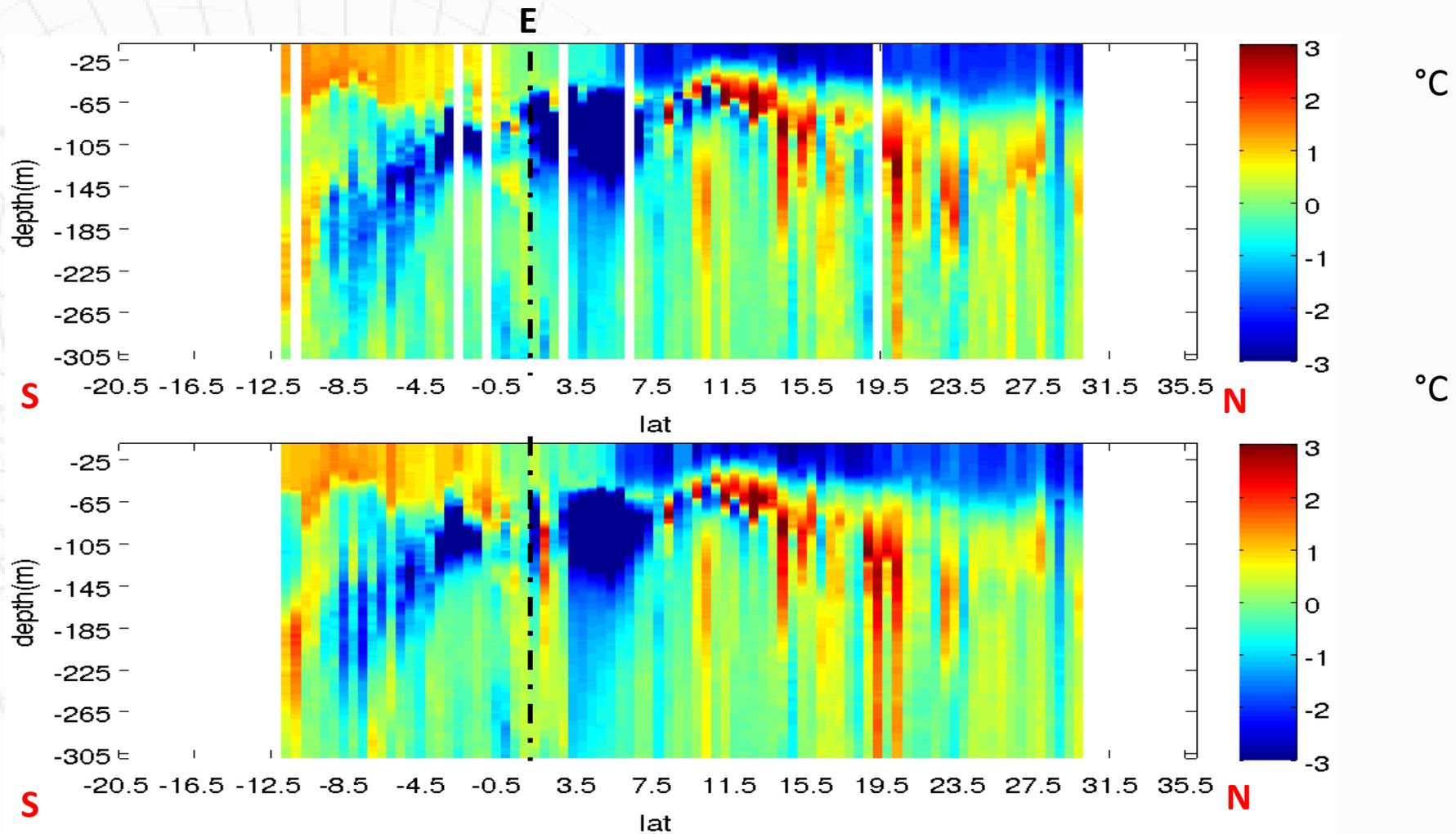
SST



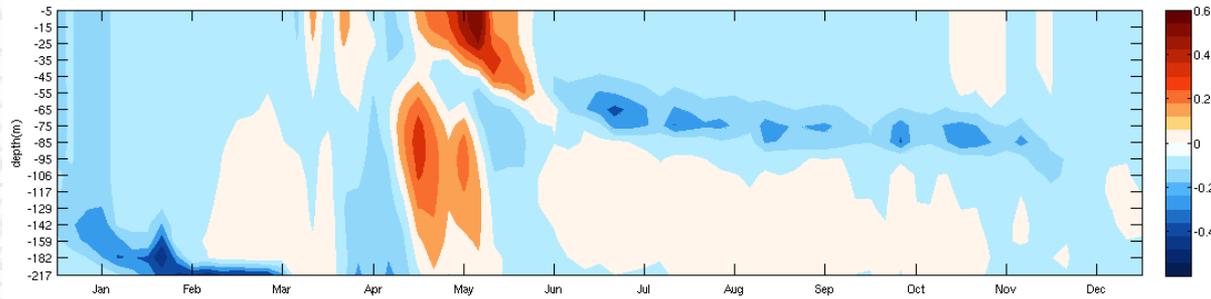
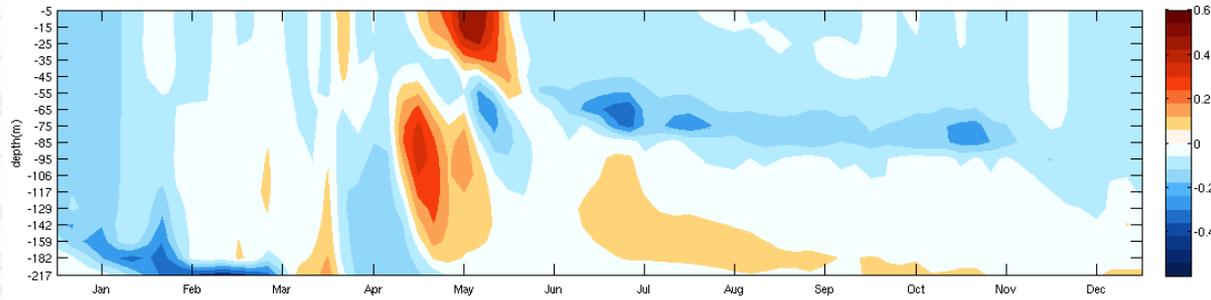
SCHL



APPLICATION: ARAMIS RECONSTRUCTION (2/4)



Interannuel

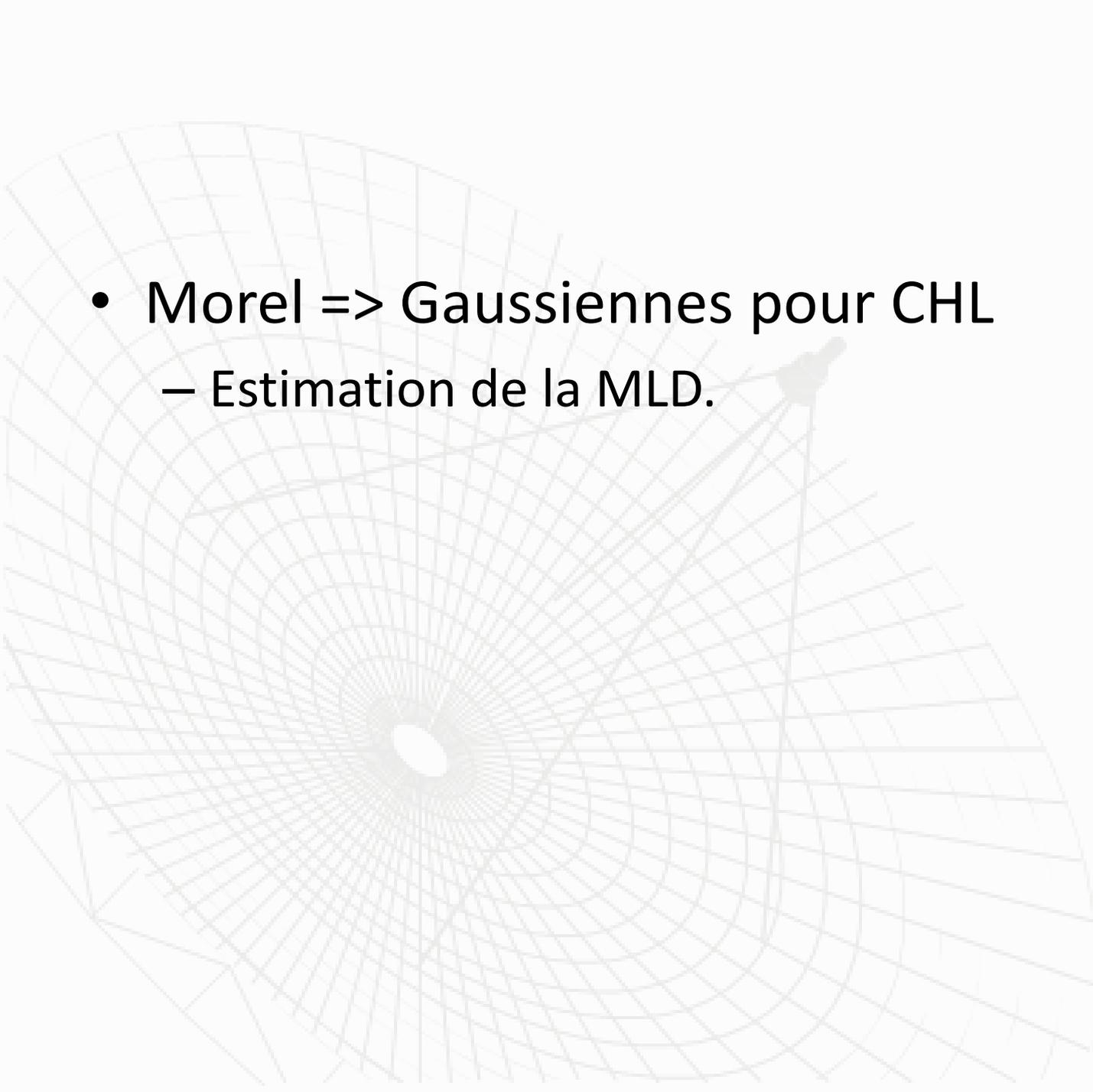


RMS ARAMIS 2

	RMSE (°C)	5 - 35m	35 - 65m	65 - 95m	95 - 125m	125 - 155m	155 - 185m	185 - 215m	215 - 245m	245 - 275m	275 - 305m	Moyen
Intervalles	Total	0,072	0,060	0,055	0,050	0,040	0,059	0,068	0,118	0,034	0,048	0,060
Moyennées	Validation	0,143	0,141	0,300	0,381	0,303	0,419	0,161	0,464	0,046	0,059	0,242
Comparaisons	Total	0,318	0,347	0,392	0,430	0,447	0,623	0,772	0,931	0,582	0,245	0,509
Point / point	Validation	0,428	0,563	0,859	0,858	0,655	0,863	1,053	1,317	0,924	0,223	0,774

Performances numériques

- 16GB RAM, 8 processeurs 1.2GH:
 - Apprentissage et optimisation complète en
 - ~7h30 pour BATS
 - ~4h45 pour ARAMIS
 - Reconstruction d'une série totale
 - Temps < 1 minute

- 
- Morel => Gaussiennes pour CHL
 - Estimation de la MLD.