

M2 - STL

Algorithmes sur les séquences en bioinformatique

Cours 2: Algorithmes de comparaison des séquences

Alessandra Carbone
Université Pierre et Marie Curie

Problème : étant données deux (ou plusieurs) séquences d'ADN ou de protéines, trouver le meilleur alignement entre elles.

On commencera avec l'alignement de paires de séquences ou le problème de l'*appariement inexacte*, c-a-d qu'on permettra que certains acides-aminés ne soient pas alignés.

Les algorithmes que l'on introduira sont basés sur la technique de la **programmation dynamique**. Pour chaque algorithme on présente:

- l'intuition qui est à la base du processus récursif
- la définition formelle du processus récursif
- sa complexité

Définition du problème et motivation biologique

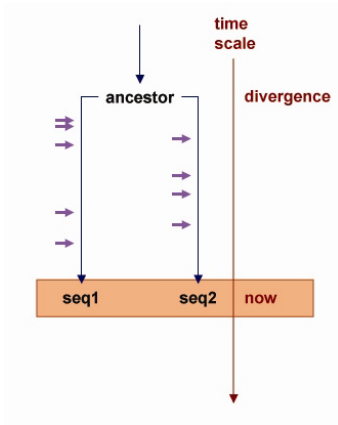
Plusieurs problèmes en bioinformatique demandent une solution du problème de l'alignement entre séquences:

- Reconstruction de séquences longues d'ADN à partir de nombreuses séquences courtes
- Comparer deux ou plusieurs séquences par similarité: recherche de séquences similaires dans une bdd
- Organiser, mémoriser, récupérer et comparer des séquences d'ADN dans les bdd
- Explorer des motifs souvent récurrents dans les séquences nucléotidiques
- Trouver des éléments informatifs dans les protéines (e.g. les domaines fonctionnels) ou les séquences d'ADN (e.g. les sites de regulation)

Homologie de séquences : similarité et différence

Hypothèse théorique. Tout le matériel génétique a un ancêtre commun. La plupart des différences entre familles d'espèces sont dues aux modifications locales entre séquences nucléotidiques:

- **insertion** – insertion d'une ou plusieurs bases dans une séquence
- **délétion** – délétion d'une ou plusieurs bases dans la séquence
- **substitution** – remplacement d'une base par une autre base



Une séquence **ACGTACGT** a pu évoluer pour donner :

-- ACG- T- A- --CG- T-----
A CACGGTCC TAA TAA TGGCC

--- AC- GTA -C --G -T --
CAG -GA AGA TCTT A GT TC

simulation: après 9500 générations avec des probabilités de

délétion: 0.0001

insertion: 0.001

substitution A/G, T/C: 0.00008

substitution A/C, T/G: 0.00002

Ces deux séquences ont le même ancêtre et on veut connaître leur alignement qui reflète leur origine commune (c'est-à-dire **ACGTACGT**)

Par exemple, l'alignement pourrait être

-AC AC- GGTCCTAA T- -AATGGCC
CAG -GAA -G-AT- -CTTAGTTC - -

Distance et similarité

Etant données deux séquences, on associe **des poids** aux paires d'acides-aminés/espaces qui représentent le coût des modifications.

La **distance** entre séquences est la **somme minimale** des poids pour un ensemble de modifications qui transforment l'une dans l'autre.

La **similarité** entre deux séquences est la **valeur maximale** de la somme des poids.

Modèles d'alignement

Problème 1: (alignement globale)

Etant données deux séquences S et T ayant à peu près la même longueur, quelle est la similarité maximale entre elles? Trouver le meilleur alignement.

Problème 2: (alignement locale)

Etant données deux séquences S et T, de longueur possiblement différente. Quelle est la similarité maximale entre une sous-séquence de S et une sous-séquence de T? Trouver les sous-séquences les plus similaires.

Problème 3: (globale-locale - alignement libre aux extrémités)

Etant données deux séquences S et T, de longueur possiblement différente. Trouver le meilleur alignement entre sous-séquences de S et de T quand au moins une de ces sous-séquences est un préfix de la séquence originale et une (pas forcément l'autre) est un suffixe.

Définition: un **gap** est l'intervalle d'espaces maximale d'une séquence donnée par rapport à un alignement. La longueur du gap est le nombre d'opérations d'indels effectuées sur la séquence.

Définition: une **fonction de pénalité de gap** est une fonction que mesure le coût d'un gap comme une fonction non linéaire de sa longueur.

Problème 4: (pénalité de gap)

Etant données deux séquences S et T, de longueur possiblement différente. Trouver le meilleur alignement entre les deux séquences utilisant la fonction de pénalité de gap.

Alignement globale optimale en $O(nm)$ en utilisant la programmation dynamique

Entrées: deux séquences S, T tels que $|S|=n$, $|T|=m$ et $n \sim m$

Sortie: un alignement optimale et sa valeur

On va résoudre un problème plus compliqué :

$V(i,j)$ = valeur de l'alignement optimale de $S[1], \dots, S[i]$ avec $T[1] \dots T[j]$ pour tout $0 \leq i \leq n$, $0 \leq j \leq m$

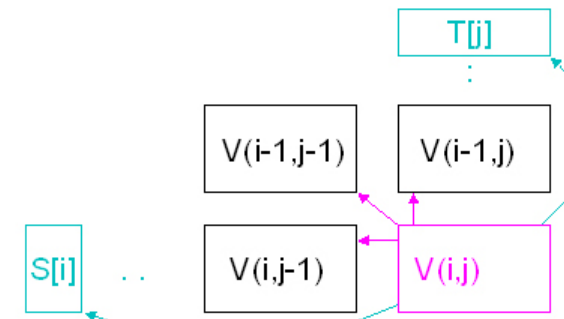
Base:

$$V(i,0) = \sum_{k=0}^i \sigma(S[k], -)$$

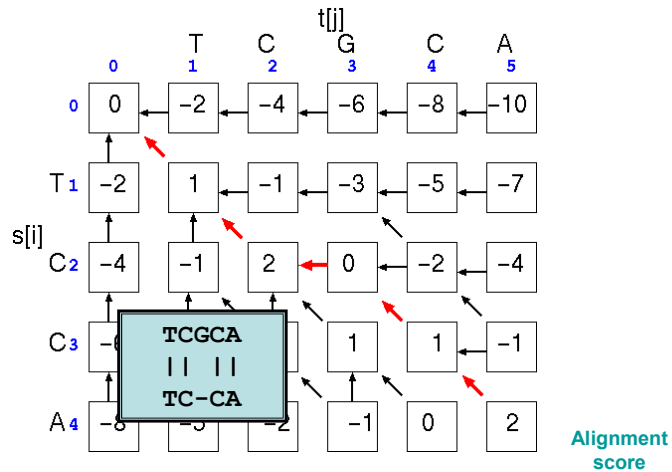
$$V(0,j) = \sum_{k=0}^j \sigma(-, T[k])$$

Recurrence:

$$V(i,j) = \max \left\{ \begin{array}{l} V(i-1,j-1) + \sigma(S[i], T[j]) \\ V(i-1,j) + \sigma(S[i], -) \\ V(i,j-1) + \sigma(-, T[j]) \end{array} \right.$$



Sequence alignment algorithm : Needleman & Wunsch



On peut prouver par induction que cette définition donne bien un alignement optimal.

Exemple: calcul sur tableau

Mismatch = -1
Match = 2

j	0	1	2	3	4	5
i		c	a	d	b	d
0	0	-1	-2	-3	-4	-5
1	a	-1	-1	1		
2	c	-2	1			
3	b	-3				
4	c	-4				
5	d	-5				
6	b	-6				

← T

Time = O(mn)

↑ S

Recherche d'un alignement par backtracking

j	0	1	2	3	4	5
i		c	a	d	b	d
0	0	-1	-2	-3	-4	-5
1	a	-1	-1	1	0	-1
2	c	-2	1	0	0	-1
3	b	-3	0	0	-1	2
4	c	-4	-1	-1	1	1
5	d	-5	-2	-2	1	0
6	b	-6	-3	-3	0	3

← T

↑ S

Théorème: Le temps de calcul de l'algorithme est $O(nm)$. L'espace de calcul est $O(n+m)$, si seulement $V(S,T)$ est demandé, si non $O(nm)$ pour la reconstruction de l'alignement.

Preuve: devoir.

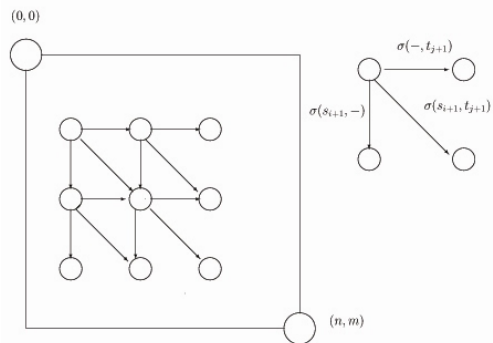
Graphe d'alignement

Il est souvent utile de représenter à l'aide d'un **graphe pondéré** les solutions obtenues par programmation dynamique de problèmes sur les séquences.

Définition: Soient S et T deux séquences de longueur n et m . Un graphe d'alignement est un graphe dirigé $G=(V,E)$ de $(n+1) \times (m+1)$ noeuds, chacun étiqueté par une paire (i,j) distincte ($0 \leq i \leq n$, $0 \leq j \leq m$), et avec arêtes pondérées définies comme suit:

1. $((i,j),(i+1,j))$ avec poids $\sigma(S[i+1], -)$
2. $((i,j),(i,j+1))$ avec poids $\sigma(-, T[j+1])$
3. $((i,j),(i+1,j+1))$ avec poids $\sigma(S[i+1], T[j+1])$

Un chemin du noeud $(0,0)$ au noeud (n,m) dans le graphe d'alignement correspond à un alignement et le poids totale est la valeur de l'alignement. Le but est de trouver le chemin le plus cher de $(0,0)$ à (n,m) .



Trois possibilités représentant chacune une étape de l'algorithme.

Alignement globale en temps linéaire

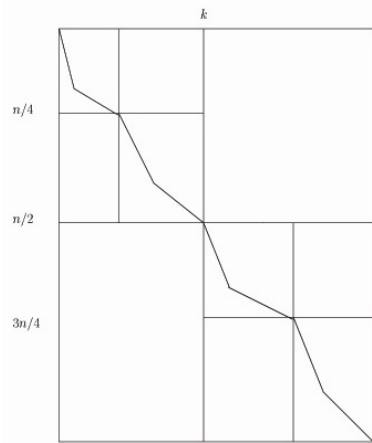
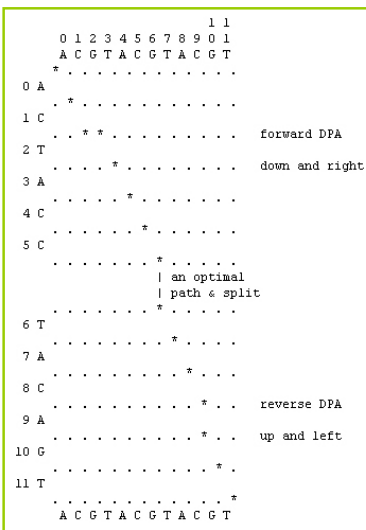
L'algorithme de backtracking appliqué à la recherche d'un chemin, demande la mémorisation de la matrice entière. En conséquence, la complexité en espace est $O(nm)$. Hirschberg a proposé une solution qui donne une complexité d'espace en $O(\min(n,m))$ basée sur la méthode « divide et impera ».

L'algorithme:

Dénoter $V^r(i,j)$ la valeur de l'alignement optimale des derniers i caractères de S contre les derniers j caractères de T .

1. Calculer $V(A,B)$ et sauver en mémoire les valeurs de la $n/2$ -ième ligne, où $V(A,B)$ dénote la *matrice en avant* F .
2. Calculer $V^r(A^*,B^r)$ et sauver en mémoire les valeurs de la $n/2$ -ième ligne, où $V^r(A,B)$ dénote la *matrice en arrière* B .
3. Trouver la colonne k^* qui croise le point $(n/2, k^*)$ satisfaisant:

$$V(n/2, k^*) + V^r(n/2, m - k^*) = V(n, m)$$
4. Partitionner récursivement le problème en deux sous-problèmes:
 - (i) trouver le chemin de $(0,0)$ à $(n/2, k^*)$
 - (ii) trouver le chemin de (n, m) à $(n/2, m - k^*)$



Lemme: $V(n,m) = \max_{0 \leq k \leq m} \{V(n/2,k) + V'(n/2,m-k)\}$

Chaque calcul de programmation dynamique demande la mémorisation de la ligne « n/2 » (la ligne du milieu), mais une fois qu'on trouve le point k* alors cette ligne peut être déchargée. Si n < m on peut mémoriser la colonne du milieu à la place de la ligne.

On en déduit que la complexité d'espace est $O(\min(n,m))$.

Alignement locale

Entrées : deux séquences S et T.

Sortie : la sous-séquence de S et la sous-séquence de T de similarité maximale.

Pourquoi rechercher en locale ?

Pour ignorer les parties non-codantes: Dans l'**ADN**, les régions non-codantes (introns) sont plus sujettes aux mutations que les régions codantes (exons). Quand on cherche un alignement locale entre deux séquences d'ADN, on trouve souvent que le meilleur match concerne des exons.

Pour détecter les domaines des protéines: **protéines** différentes peuvent montrer des similarités locales, les "homeoboxes". Il s'agit dans la plus part des cas de sous-unités fonctionnelles d'une protéine.

Problème de l'alignement locale de suffixes: étant données deux séquences S et T et deux indices i et j, trouver un suffixe (possiblement vide) de $S[1\dots i]$ et un suffixe (possiblement vide) de $T[1\dots j]$ tels que la valeur de leur alignement est maximale sur tous les alignements de suffixes de $S[1\dots i]$ et $T[1\dots j]$.

La solution du problème de l'alignement locale revient à trouver la solution maximale du problème de l'alignement locale des suffixes sur tous les indices i et j de S et T.

Terminologie et restriction:

$V(i,j)$ dénote la valeur de l'alignement local des suffixes optimale pour la paire d'indices i, j

Les poids des opérations d'édition sont tels que:

$$\sigma(x,y) = \begin{cases} \geq 0 & \text{si } x,y \text{ font un match} \\ \leq 0 & \text{si } x,y \text{ ne font pas de match ou si l'un des deux est un espace} \end{cases}$$

Schema de l'algorithme

1. Calcule l'alignement des suffixes locale (pour tout i et j) de S[1...i] et T[1...j]. Cela est fait en utilisant l'algorithme d'alignement globale.
2. Recherche les résultats et trouve les indices i* et j* de S et T, après lesquels la similarité peut seulement décroître.

Définition récursive

Base pour tout i,j $V(i,0)=0, V(0,j)=0$

Recurrence
$$V(i,j) = \max \begin{cases} 0 \\ V(i-1,j-1) + \sigma(S[1\dots i], T[1\dots j]) \\ V(i-1,j) + \sigma(S[1\dots i], -) \\ V(i,j-1) + \sigma(-, T[1\dots j]) \end{cases}$$

Calcule i* et j* : $V(i^*, j^*) = \max_{1 \leq i \leq n, 1 \leq j \leq m} V(i,j)$

A.Carbone-UPMC

29

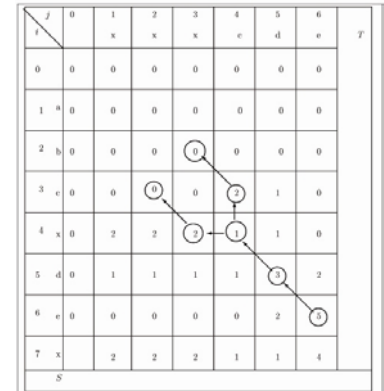
Le 0 dans la base et dans la récurrence permet à l'alignement de commencer n'importe où dans les séquences, plutôt que le forcer à commencer au début des séquences. La recherche de i* et j* qui donne le meilleur $V(i^*, j^*)$ permet à l'alignement de terminer n'importe où dans les séquences, plutôt que de terminer à la fin des séquences. Ces deux conditions prises ensemble permettent d'obtenir le meilleur alignement de deux sous-séquences de S et T.

Exemple: $\sigma(\text{match}) = 2$
 $\sigma(\text{mismatch}) = -1$

Meilleurs sous-séquences:

C- DE
 CXDE

XCDE
 X-DE



A.Carbone-UPMC

30

L'alignement locale peut être résolu en :

temps $O(nm)$: le nombre constante d'opérations pas cellule du tableau pour le calcul de $V(i,j)$, donne un temps $O(nm)$ pour le remplissage du tableau. La recherche de $V(i^*, j^*)$ demande un temps $O(nm)$.

espace linéaire $O(n + m)$: avec l'algorithme de Hirschberg.

A.Carbone-UPMC

31

Alignement libre aux extrémités

Entrées: deux séquences S et T, de longueur possiblement différente.
Sortie: le meilleur alignement entre sous-séquences de S et de T quand au moins une de ces sous-séquences est un préfix de la séquence originale et une (pas forcément l'autre) est un suffixe.

Tout opération d'indel aux extrémités des séquences a coût 0. Cette condition permet d'oublier le début et la fin des séquences dans l'alignement et elle permet l'alignement de séquences qui se chevauchent ou quand l'une inclue l'autre.

Exemple: soit S= CACTGTAC et T= GACACTTG
 Supposons que $\sigma(\text{match})=2$ et que $\sigma(\text{mismatch}) = -1$.

L'alignement globale optimale a valeur 1 : L'alignement locale optimale a valeur 9:

CAC - - T - GTAC
 |||
 GACTT TG - - -

- - CAC- TGTAC
 |||
 GACACTTG - - -

A.Carbone-UPMC

32

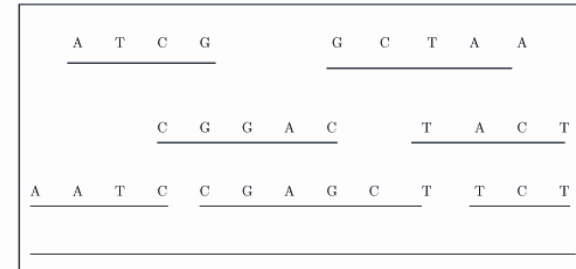
Motivations:

Dans le problème de l'**assemblage de séquences shotgun** nous avons un grand nombre de séquences qui se chevauchent partiellement et qui proviennent de plusieurs copies de la même séquence d'ADN qui est inconnue et que l'on veut reconstruire. Le problème revient à utiliser la comparaison de sous-séquences pour déduire la séquence d'origine (en utilisant la partie chevauchante pour coller ensemble les différents bouts).

Deux sous-séquences qui proviennent de deux parties différentes de la séquence d'origine auront une valeur d'alignement globale très basse, ainsi que la valeur de l'alignement local.

Comme deux sous-séquences chevauchantes n'ont pas, avec une forte probabilité, la même position initiale (ni la même position finale) sur la séquence originale, elles auront une valeur d'alignement globale qui sera toujours basse. Par contre, l'alignement libre aux extrémités **sera élevé** parce qu'elles possèdent une partie chevauchante. Le chevauchement sera détecté, et les sous-séquences seront collées ensemble en utilisant l'alignement trouvé.

De façon similaire on traite le cas d'une séquence qui contient une autre séquence.



A.Carbonne-UPMC

34

Algorithme d'alignement à extrémités libres

On fixe les conditions initiales pour permettre poids 0 **aux extrémités d'au plus une séquence**.

Après avoir rempli la table avec les valeurs $V(i,j)$ on recherche la valeur maximale dans la ligne plus en bas/colonne plus à droite, permettant de terminer au plus une séquence avant l'autre, avec poids zéro pour toutes les opérations d'insertion de gaps de là jusqu'à la fin. La valeur obtenue est la valeur optimale.

La séquence alignée est déterminée à partir de la case (0,0) jusqu'à la fin de la table (ligne plus en bas/colonne plus à droite). Après ce point, toute opération d'introduction de gap jusqu'à (n,m) ne sont pas tenues en compte dans la valeur totale (même si elle sont présentes dans la table).

A.Carbonne-UPMC

35

Base	pour tout i, j $V(i,0)=0, V(0,j)=0$
Réurrence	$V(i,j) = \max \begin{cases} V(i-1,j-1) + \sigma(S[1\dots i], T[1\dots j]) \\ V(i-1,j) + \sigma(S[1\dots i], -) \\ V(i,j-1) + \sigma(-, T[1\dots j]) \end{cases}$
Calcul de i^* :	$V(i^*,m) = \max_{1 \leq i \leq n, m} V(i,j)$
Calcul de j^* :	$V(n,j^*) = \max_{n, 1 \leq j \leq m} V(i,j)$
Définition du score d'alignement:	$V(S,T) = \max \{V(n,j^*), V(i^*,m)\}$

Complexité en temps:

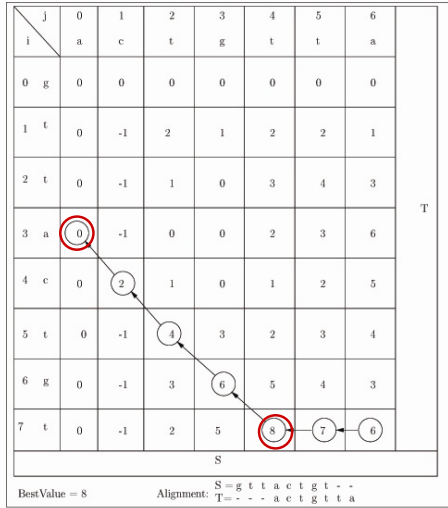
trouver la matrice $O(nm)$
trouver i^* et j^* $O(n+m)$

Complexité en espace:

calcul de la matrice $O(n+m)$ avec la méthode de Hirschberg
calcul de i^* et j^* $O(n+m)$ pour le stockage de la dernière ligne/colonne

A.Carbonne-UPMC

36



Pénalités de gaps

Jusqu'à maintenant les éléments qui ont permis la définition du score ont été matches, mismatches et espaces. Un élément important qui est souvent pris en compte est le « gap », c'est-à-dire **une suite consécutive maximale d'espaces dans un alignement**.

L'idée est de traiter le gap comme un élément unique plutôt que de lui associer une valeur qui est la somme d'espaces. Il y a plusieurs façons de traiter la valeur d'un gap et on présentera plusieurs modèles.

#gaps = nombre de gaps dans un alignement

Exemple: l'alignement suivant à 8 gaps et un total de 11 espaces.

- ACAC- GGCCTA A T-- AATGGCC
 | | | | | | | |
 CAG -GAA -G-A T-- CTTAGTTC --

Motivations

L'insertion ou la délétion de sous-séquences entières a lieu souvent comme un **seul événement mutational**. En plus, plusieurs de ces événements mutationnels peuvent créer des gaps de tailles différentes. Pour éviter d'associer des valeurs qui sont trop importantes à ces événements mutationnels il faut que l'on traite les gaps comme des unités en soit même.

Exemple: considérons l'alignement d'ARN messagers (qui codifient une protéine, sans contenir des traces d'introns) avec des séquences d'ADN qui contiennent exons et introns. On sait que l'alignement contiendra plusieurs gaps possiblement très longs correspondant aux introns. L'utilisation d'**une bonne pénalité de gap nous permettra d'éviter des scores très faibles** pour ces alignements et donc nous permettra de trouver l'alignement vrai.

Modèle à pénalité de gap constante : espace n'a pas de poids
 gap a poids W_g indépendant de sa longueur

Le problème revient à maximiser $\sum \sigma(S'[i], T'[i]) + W_g \times \text{\#gaps}$

où S' et T' sont les séquences S et T avec espaces.

Modèle à pénalité de gap affine : le gap a 2 poids
 W_g ouverture du gap
 W_s extension du gap avec un espace
 avec un poids totale $W_{\text{totale}} = W_g + qW_s$
 pour un gap de longueur q

Le problème revient à maximiser $\sum \sigma(S'[i], T'[i]) + W_g \times \text{\#gaps} + W_s \times \text{\#espaces}$

Algorithme avec pénalité de gap affine

Pour aligner S et T, on considère les préfixes S[1...i] et T[1...j]. Un alignement des deux préfixes rentre dans l'une des catégories suivantes:

1. S ----- i
T ----- j

où les caractères S[i] et T[j] sont alignés entre eux.
Cette hypothèse inclut les cas S[i]=T[j] et S[i] ≠ T[j].

2. S ----- i _____
T ----- j

où le caractère S[i] est aligné à un caractère strictement à la gauche de T[j].
L'alignement termine alors avec un gap sur la droite dans S.

3. S ----- i
T ----- j _____

où le caractère T[j] est aligné à un caractère strictement à la gauche de S[i].
L'alignement termine alors avec un gap sur la droite dans T.

Notation:

G(i,j) la valeur maximale des alignements de type 1
E(i,j) la valeur maximale des alignements de type 2
F(i,j) la valeur maximale des alignements de type 3
V(i,j) la valeur maximale de l'alignement

Comme dans les autres algorithmes d'alignement on donne une définition récursive de la table d'alignement. Il y a 3 relations de récurrence à définir, pour G(i,j), E(i,j) et F(i,j).

Exemple: on considère E(i,j). On doit regarder les opérations d'insertion d'un espace et assigner la bonne valeur : cela dépend du poids des espaces (qW_s), mais aussi du poids de l'ouverture d'un gap (W_g). Elle sera calculée en regardant les valeurs de la table calculées aux récursions précédentes. Il y a deux cas possibles:

1. S termine à la gauche de T, comme pour l'étape courante. Nous ajoutons un poids d'extension à la valeur précédente: $E(i,j-1)+W_s$.
2. S et T terminent au même endroit (type 1). Nous ajoutons un poids d'ouverture et un poids d'extension: $V(i,j-1)+W_g+W_s$.

Définition récursive

Base

$$\begin{aligned} V(0,0) &= 0 \\ V(i,0) &= E(i,0) = W_g + iW_s \\ V(0,j) &= F(0,j) = W_g + jW_s \end{aligned}$$

Réurrence

$$\begin{aligned} V(i,j) &= \max \{ E(i,j), F(i,j), G(i,j) \} \\ G(i,j) &= V(i-1, j-1) + \sigma(S[i], T[j]) \\ E(i,j) &= \max \{ E(i, j-1) + W_s, V(i, j-1) + W_g + W_s \} \\ F(i,j) &= \max \{ F(i-1, j) + W_s, V(i-1, j) + W_g + W_s \} \end{aligned}$$

Complexité en temps: $O(nm)$ parce qu'on calcule 4 matrices à la place d'une seule.

Complexité en espace: $O(nm)$ pour une implémentation de base.

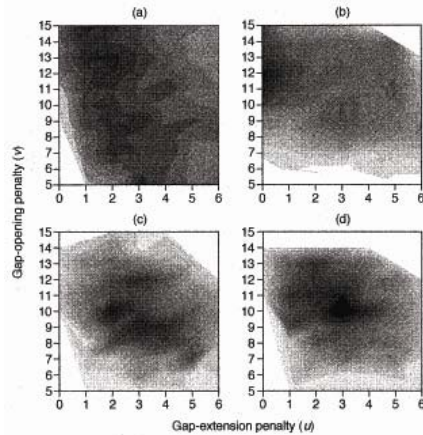
Modèle à pénalité de gap convexe : chaque espace ajouté à un gap contribue moins au poids que le précédent. On utilise par exemple la fonction de poids $W_g \cdot \log(q)$, où q est la longueur du gap.

Le problème est résoluble en temps $O(nm \log(m))$.

Modèle à pénalité de gap arbitraire : n'importe quelle fonction de poids est acceptable, et le poids est fonction de la longueur q du gap.

Le problème est résoluble en temps $O(nm(m+n))$.

Distribution de meilleures valeurs de gap pour des programmes d'alignement différents



- (a) Pairwise Alignment
- (b) ClustalW
- (c) Randomized Iterative Method
- (d) Double Nested Iterative Method

Profils calculés sur 54 familles, ou les méthodes ont pu rejoindre le 99% de paires correctement alignées.

A.Carbone-UPMC

En 3eme dimension: pourcentage de paires correctement alignées

45

Why « pairwise alignment » again ?

Detection of homologous proteins

Genomes with the 60% of genes of unknown function
Ex: *Plasmodium falciparum*

Detection of protein binding sites

Larger families of homologous proteins allow to study a broader variety of evolutionary signals.

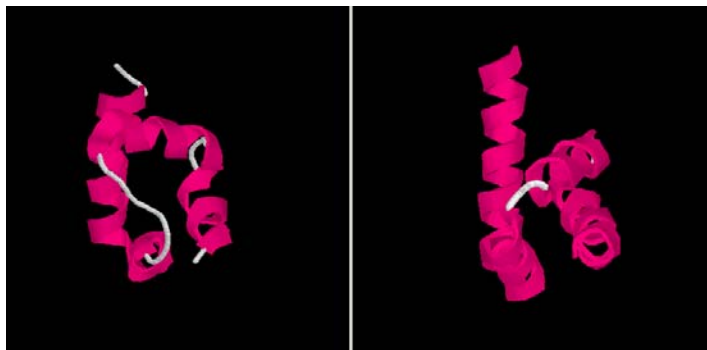
Sensitive multiple alignment

Good pairwise alignment is demanded to avoid error propagation especially for multiple alignment of sequences with low homogy

A.Carbone-UPMC We look for very weak signals

46

An example: looking for proteins with similar structure

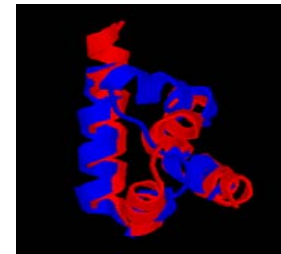


protéine ribosomale L20

A.Carbone-UPMC Aquifex aeolicus

protéine polyA binding Homo sapiens₄₇

Superposition of the two structures



bacteria and human

Structural alignment program PROSUP

```
bacteria: |---WIARINAAVRA--YGLNYSTFINGLKKAGIELDRKILADMAVRDPQAFEQVNVKKEALQVQ-
equiv.: | *****
human: |HRQALGERLYPRVQAMQPAFASKITG-----MLLELSPAQLLLLASEDSLRRVDEAMELI IAHG
Identity: | # # # # #
```

A.Carbone-UPMC

Number of Identities = 5; RMSD = 2.41Å

48

Sequence alignment programs fail

Sequence alignment program CLUSTALW

```
bacteria: |-WIARINAAVRAYGLNYSTFINGLKKAGIELDRKILADMVARDPOAFEQVNVKVKALQVQ
human:  |HRQALGERLYPRVQAMQPAFASKITGMLELLESPAQLLLLLSEDSLRARVDEAMELIIAAG
Identity: | # # # # # #
```

Number of identities = 6

A.Carbone-UPMC

49

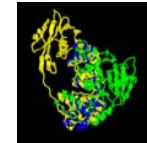
Alignment and Homology Search Reliability

A)

Structure



Good reliability

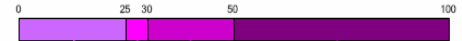


B)

Sequence

MFPDAHCELVHRNFPELLIAVVLSAQ

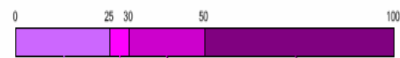
MFPDAHCEL--VHRNFPELLIAVVLSAQ
PGLRLPGC----VDAFEGVRAILGQL
YHDNEGWFPETDSKILFEMI CLEGGQAG



- 0-25%
 - Homology possible but often unclear
 - Detection and Alignment difficult
- 25-30%
 - Highly likely homology
 - Detection and Alignment becoming tricky
- 30-50%
 - Good homology
 - Relatively easy detection and alignment
- >50%
 - Close homology
 - Trivial detection and alignment

A.Carbone-UPMC

E.E. Hill and S. E. Brenner, 50
IPAM Structural Proteomics, 2004



- 0-25%
 - Homology possible but often unclear
 - Detection and Alignment difficult
- 25-30%
 - Highly likely homology
 - Detection and Alignment becoming tricky
- 30-50%
 - Good homology
 - Relatively easy detection and alignment
- >50%
 - Close homology
 - Trivial detection and alignment

Introduction of structural information in

1) sequences

2) alignment parameters:

- substitution matrices
(Overington et al., 1992, Teodorescu et al., 2004)
- gap penalties
(Lesk et al., 1986)

A.Carbone-UPMC

51

```

10 20 30 40 50 60 70
lh88_C2/1-73  -LIGPWKEEDRVIELVQRY  SPNR  -NSVIATH  LKGIQGCCBWWHNLNPE  -----
lmse-1/1-73   -MLINGPWKEEDRVIELVQRY  SPNR  -NSVIATH  LKGIQGCCBWWHNLNPEVK  -----
lgv2-1/1-73   -ELIIGPWKEEDRVIELVQRY  SPNR  -NSVIATH  LKGIQGCCBWWHNLNPEVK  -----
lgv5/1-73     -LIGPWKEEDRVIELVQRY  SPNR  -NSVIATH  LKGIQGCCBWWHNLNPE  -----
lgv4/1-73     -LIGPWKEEDRVIELVQRY  SPNR  -NSVIATH  LKGIQGCCBWWHNLNPE  -----
lmbg/1-73     -LIGPWKEEDRVIELVQRY  SPNR  -NSVIATH  LKGIQGCCBWWHNLNPE  -----
lmbh/1-73     -LIGPWKEEDRVIELVQRY  SPNR  -NSVIATH  LKGIQGCCBWWHNLNPE  -----
lh8a_C1/1-73  -PELNGPWKEEDRVIELVQRY  SPNR  -NSVIATH  LKGIQGCCBWWHNLNPEVK  -----
la5j-1/1-73  -GIDLVGPWKEEDRVIELVQRY  QTQQ  -WTLIATH  LKGIQGCCBWWHNLNPEVK  -----
lguu/1-73     -LGIHTRWREEDRKLKLLVEGN  GTTD  -KVVIAMN  LPNRIDVQCHHWAKVLNPE  -----
lh88_C1/1-73  -MGLHGTWREEDRKLKLLVEGN  GTTD  -KVVIAMN  LPNRIDVQCHHWAKVLNPE  -----
lmbf/1-73     -LGIHTRWREEDRKLKLLVEGN  GTTD  -KVVIAMN  LPNRIDVQCHHWAKVLNPE  -----
liiy/1-73     -RARRQAWLWEEDKNLRSVFRFY  GEGN  -KSKILLHYEFNNRQSVLIDPWRTWKRLIISDSD  -----
liiv/1-73     -RARRQAWLWEEDKNLRSVFRFY  GEGN  -KSKILLHYEFNNRQSVLIDPWRTWKRLIISDSD  -----
lba5/1-73     -RRQAWLWEEDKNLRSVFRFY  GEGN  -KSKILLHYEFNNRQSVLIDPWRTWKRLIISDSD  -----
lgv2-2/1-73  -TSWFEEDRIIQAARHLL  GNN  -NAEIALL  LPKIDNAVNHNSMFRFVY  -----
lmbk/1-73     -KSTWFEEDRIIQAARHLL  GNN  -NAEIALL  LPKIDNAVNHNSMFRFVY  -----
lh88_C3/1-73  -KSTWFEEDRIIQAARHLL  GNN  -NAEIALL  LPKIDNAVNHNSMFRFVY  -----
lmbj/1-73     -KSTWFEEDRIIQAARHLL  GNN  -NAEIALL  LPKIDNAVNHNSMFRFVY  -----
lmse-2/1-73  -TSWFEEDRIIQAARHLL  GNN  -NAEIALL  LPKIDNAVNHNSMFRFVY  -----
lihy/1-73     -MEVXKTSWFEEDRIIQAARHLL  GNN  -NAEIALL  LPKIDNAVNHNSMFRFVY  -----
lh8a_C2/1-73  -TSWFEEDRIIQAARHLL  GNN  -NAEIALL  LPKIDNAVNHNSMFRFVY  -----
la5j-2/1-73  -SSWFEEDRIIFEAEVLL  GNN  -NAEIALL  LPKIDNAVNHNSMFRFVDT  -----
liem/1-73     -ESHNASFDEEDDFILDVYRLN  PTKRTHLLDDEL  VPHHGNSIDHFRVLYSK  -----
consensus/1-73 -PWFEEDRIIEAVKRY  GKNY  -NEKIAME  LPKITEKCCBWWHNL  -----

```

Physico-chemical properties of amino-acids:

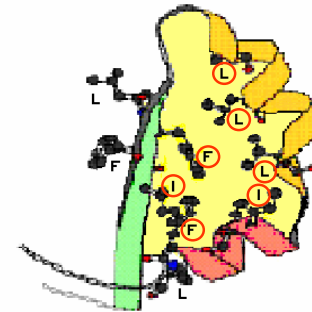
ED	ionizable	TQS	H-bonding
ILVAWMCF	non-polar/hydrophobic	HY	aromatic
RK	polar/hydrophilic	P	cyclic/proline

52

For highly divergent proteins :

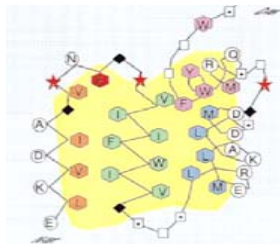
Hydrophobic amino-acids are the only ones conserved in protein families after strong evolutionary pressure

Hydrophobic Blocks



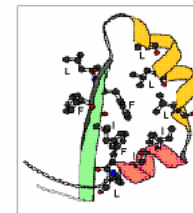
Hydrophobic core ⇨ Folding stability

Surface ⇨ Protein interactions



E K V D I A G I P N G P V I I F W I V G S T S I R E M L A K L D S D N G P T S C R W Y F T T W S Q

Prediction of Hydrophobic Blocks in Sequences



N-ter → C-ter
Specific periodicity : +/- 1, 2, 3, 4

GNATAI**IFFL**PDEGR**LQHL**ENE**L**L**THD**I**ITR**P**LE**NE**DRR**

↓

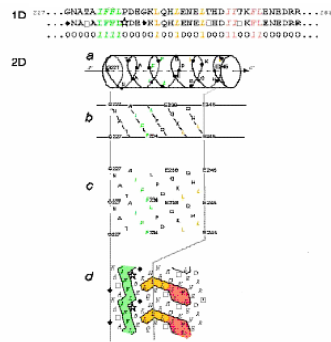
GNATA**IFFL**PDEGR**LQHL**ENE**L**L**THD**I**ITR**P**LE**NE**DRR**

Specific periodicity in sequences :

1. α -helices on protein surfaces: contacts at distance 3-4
2. β -sheets on protein surfaces: contacts at distance 2
3. Secondary structures in the hydrophobic core display chains of aa at

Hydrophobic Cluster Analysis (Mornon et al.)

α-helical 2D representation



+ manual alignment

Test of several evolutionary hypothesis on protein families with low homology

145 protein families : 613 sequences (< 30 % pairwise identity)
from the HOMSTRAD database ;
structural alignment with FUGUE of 1426 protein pairs

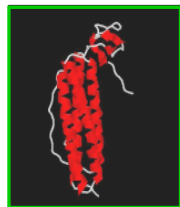
Statistics of hydrophobic blocks
Substitution matrices
Gap treatment inside and outside blocks
An alignment method based on hydrophobic blocks

Statistics on Hydrophobic Blocks

% RSS overlapping HB :
89.8%

% HB overlapping RSS :
85.7% ⇒ THC

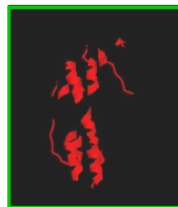
Ferritin (1RYT)



Regular Secondary Structures (RSS)

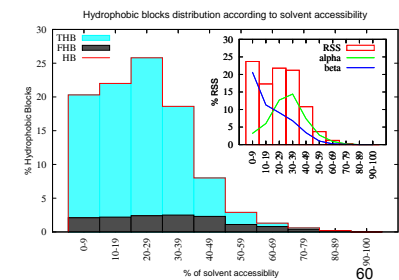
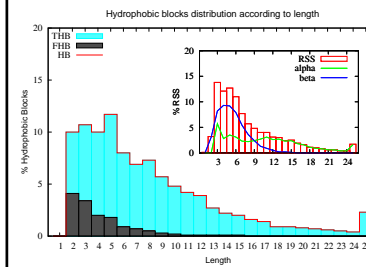


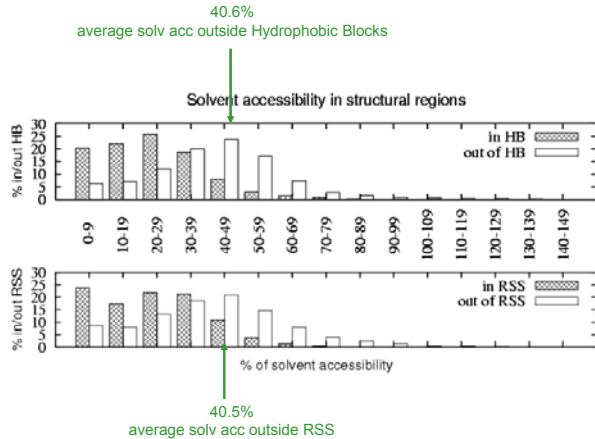
Hydrophobic Blocks overlapping RSS



Hydrophobic blocks

	RSS	HB	THB	FHB
Mean length :	8.4	8.1	8.8	4.2
Mean % solvent accessible surface	25.3	24.4	23.7	33.2





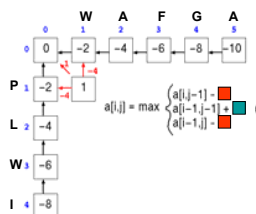
How to introduce the hydrophobic blocks signal predicted from sequences in the alignment algorithm ?

Alignment of sequences using Hydrophobic Blocks

Needleman & Wunsch

All residues are under the same evolutionary pressure

Substitution Matrix
Gap penalties

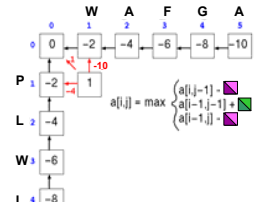


PHYBAL

Evolutionary pressure
inHBlocks > outHBlocks

inHBlocks Substitution Matrix
Gap penalties

outHBlocks Substitution Matrix
Gap penalties

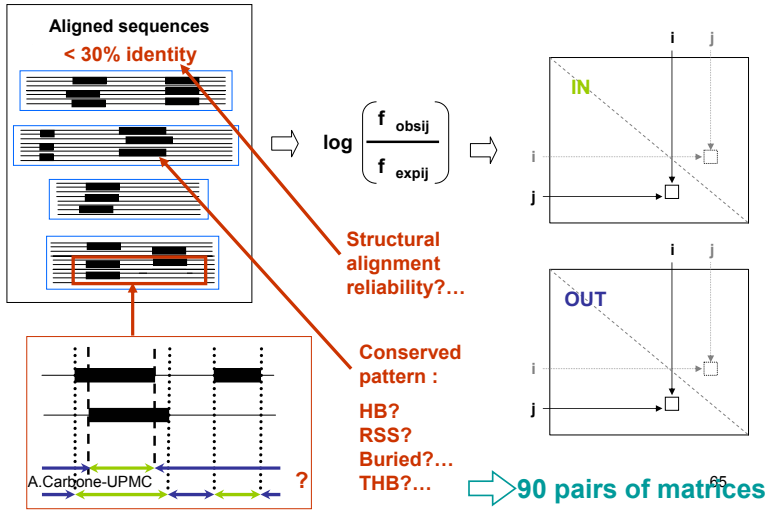


Substitution matrices construction and evaluation of the system

How to construct matrices **specific** to the hydrophobic block context ?

How to evaluate the system to observe alignment accuracy **improvement** ?

Substitution matrices construction hypothesis



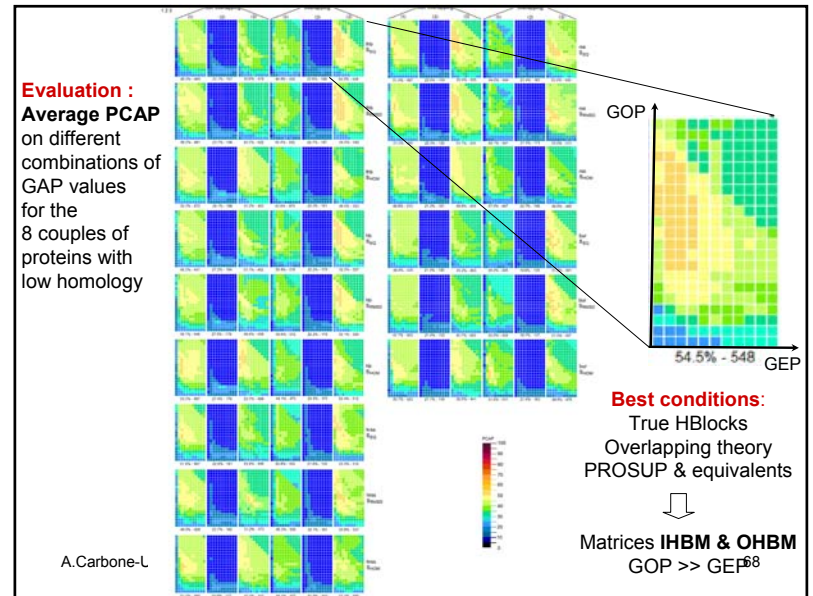
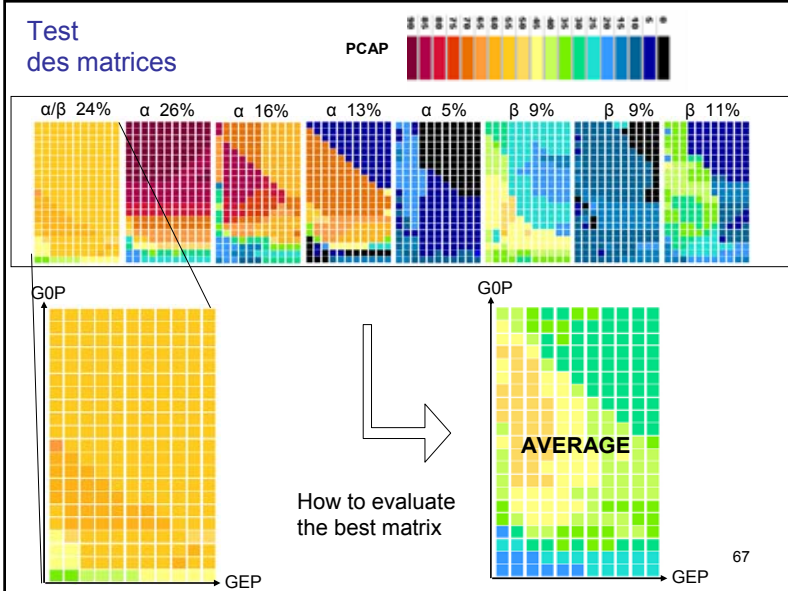
90 matrices coming from...

- A) 3 structural alignments
 - FUGUE sequence information+ substitution matrix based on structural information
 - PROSUP, smallest RMSD
 - PROSUP, largest number of equivalents, i.e. residues which are <5Å apart
 - B) 5 structures:
 - RSS
 - HB
 - buried residues (<30% solv acc)
 - THB
 - overlapping THB and RSS
 - C) 2 ways to overlap structural regions
 - D) 3 equations to define scores of substitution

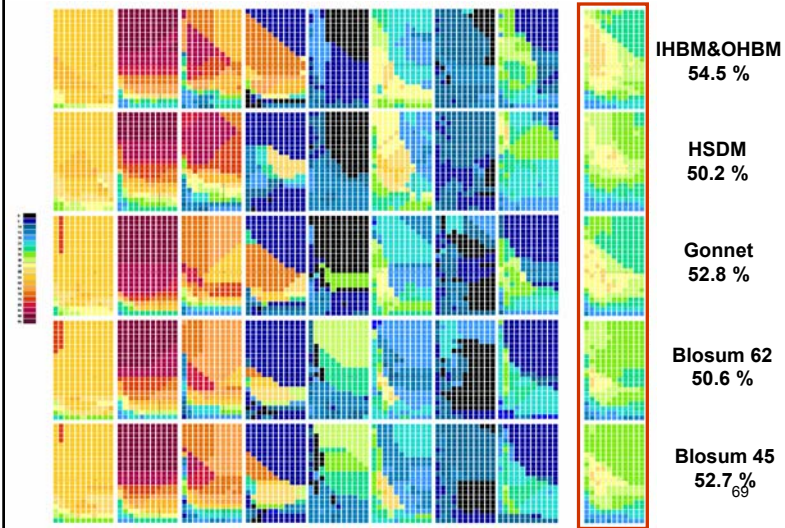
$$\log \left(\frac{f^{\text{pat}}_{\text{obs}ij}}{f^{\text{seq}}_{\text{exp}ij}} \right)$$

...and 2 more
- A. Carbone-UPMC
- 66

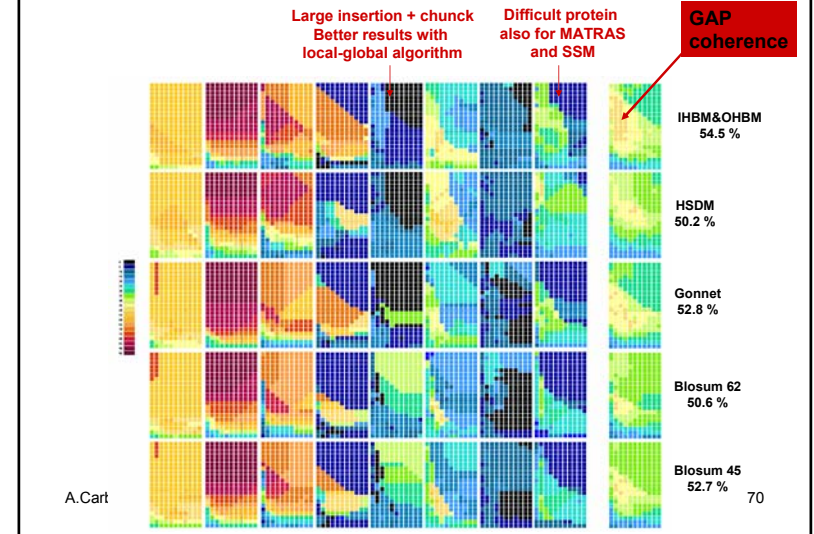
Test des matrices



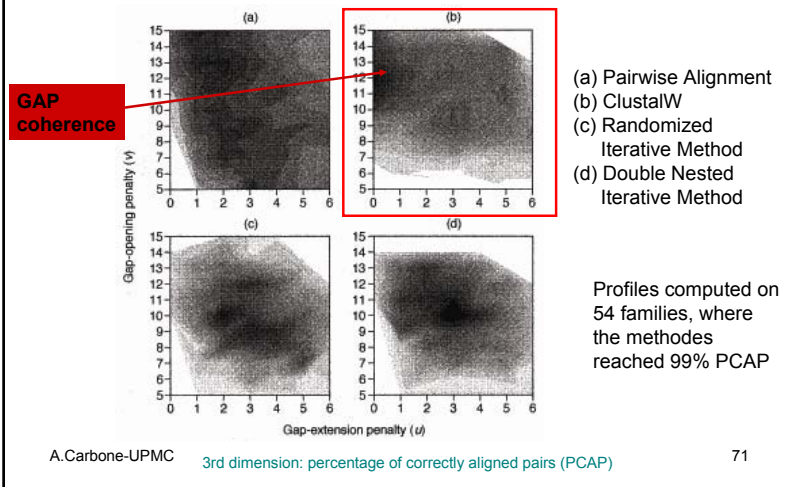
Comparison to other substitution matrices



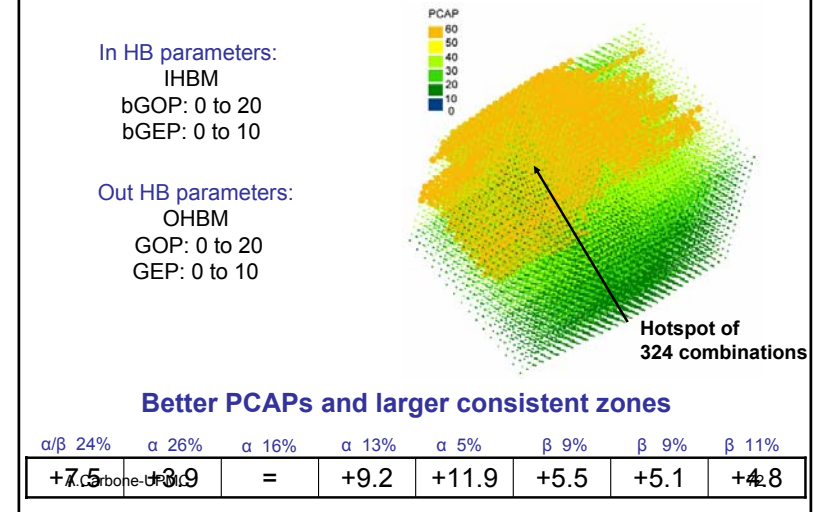
Landscape of PCAP according to gap penalties



Best gap values distribution for different alignment programs (Osamu Gotoh, JMB, 1996)



Searching for optimized gaps in 4 dimensions



Comparison with Needleman&Wunsch

PHYBAL

IHBM and OHBM

GOP/GEP : 0 to 20/10
bGOP = GOP and bGEP = GEP

2 dimensional analysis : 231 combinations
4 dimensional analysis : 324 combinations
on hotspot

Needlman & Wunsch

HSDM

Blosum30

Blosum62

Gonnet

GOP/GEP : 0 to 20/10
(231 combinations)

Global stability : (comparison with other systems) number of PCAPs which are at most 5% away from the best PCAP obtained on all systems

Local stability : (robustness of a system) number of PCAPs which are at most 2% away from the highest of the system

Domingues set

(Domingues F.S *et al.*, 2000)

165 proteins

127 protein pairs with <30% sequence identity and with a reference alignments performed by ProSup :

- at least 35 equivalent residues (i.e. residues which are <5A apart)
- share all secondary structural elements in the hydrophobic core

A.Carbone-UPMC

Method	PCAP (CAP)	glo	loc
PHYBAL hot-spot	43.7 (6461)	324	110
PHYBAL-2D + IHBM & OHBM	42.6 (6268)	49	21
Gonnet	40.2 (5982)	17	22
HSDM	43.5 (6454)	76	21
Blosum 62	41.4 (6172)	27	19
Blosum 45	41.6 (6259)	34	23
Blosum 30	38.1 (5694)	0	18
PAM250	39.3 (5797)	2	17
Johnson	41.5 (6051)	42	32
Remote Homo	32.5 (4783)	0	61

BAIiBASE (Thompson J. *et al.*, 1999)

181 proteins pairs with <25% sequence identity and ranked in 6 classes of % identity

Comparison of predicted alignments is carried on **core blocks** only, i.e. regions that can be reliably aligned.

A.Carbone-UPMC

Method	< 12	glo	loc	12-15	glo	loc	15-17	glo	loc
PHYBAL hot-spot	43.2 (296)	10	3	39.0 (1203)	11	3	40.1 (1289)	80	21
PHYBAL-2D + IHBM & OHBM	41.6 (244)	9	4	37.3 (1018)	5	4	37.7 (1211)	9	4
Gonnet	37.8 (246)	0	3	29.8 (816)	0	20	35.5 (1145)	8	2
HSDM	24.4 (176)	0	4	28.2 (823)	0	7	29.5 (1100)	0	5
Blosum 62	33.0 (239)	0	1	31.2 (964)	0	6	34.0 (1260)	0	12
Blosum 45	36.3 (243)	0	7	31.4 (898)	0	3	33.7 (1207)	0	9
Blosum 30	26.8 (197)	0	3	28.4 (830)	0	4	32.0 (1179)	0	5
PAM250	22.4 (174)	0	9	26.3 (821)	0	20	27.4 (905)	0	8
Johnson	24.8 (148)	0	14	27.8 (809)	0	7	29.1 (927)	0	29
Remote Homo	36.1 (233)	0	1	23.4 (658)	0	5	31.9 (965)	0	8

BAIiBASE (Thompson J. *et al.*, 1999)

A.Cart

	17-20	glo	loc	20-22	glo	loc	22-25	glo	loc
PHYBAL hot-spot	52.8 (2189)	161	41	57.9 (1054)	0	80	75.0 (3737)	120	42
PHYBAL-2D + IHBM & OHBM	49.4 (1971)	11	13	57.0 (1048)	0	16	72.1 (3614)	10	16
Gonnet	51.9 (2025)	10	6	55.5 (1094)	0	2	73.3 (3750)	19	16
HSDM	48.6 (1934)	1	5	54.8 (986)	0	10	75.9 (3813)	21	3
Blosum 62	49.8 (1993)	7	7	63.4 (1279)	2	1	74.9 (3852)	18	10
Blosum 45	49.1 (1917)	8	10	57.1 (1087)	0	6	72.6 (3702)	8	10
Blosum 30	45.3 (1884)	0	17	54.5 (856)	0	2	71.7 (3676)	7	11
PAM250	42.4 (1629)	0	3	54.0 (1101)	0	2	71.0 (3604)	1	8
Johnson	46.4 (1916)	0	1	55.8 (1092)	0	2	72.0 (3670)	6	8
Remote Homo	38.5 (1626)	0	7	48.7 (1042)	0	4	66.2 (3372)	0	14

76

Matrices validation : 2 matrices vs 1 matrix in 4-dim gap space

Method		≤ 12	12-15	15-17	17-20	20-22	22-25
PHYBAL	PCAP	45.0	39.3	40.1	53.3	59.9	75.1
	CAP	277	1100	1289	2164	1084	3755
	glo	570	477	1214	1003	0	282
	loc	32	102	327	185	107	132
PHYBAL hot-spot	PCAP	43.2	39.0	40.1	52.8	57.9	75.0
	CAP	296	1203	1289	2189	1054	3737
	glo	10	11	150	161	0	61
	loc	3	3	21	4	80	42
PHYBAL+Blosum62&Blosum62	PCAP	33.8	33.2	37.7	52.5	65.7	77.6
	CAP	228	1032	1258	2095	1269	3901
	glo	0	0	98	1030	58	1808
	loc	67	273	26	40	3	114

A.C.

7

Two main results come out of this analysis:

2 matrices describing hydrophobic blocks evolution

4 dimensionality of the gap space

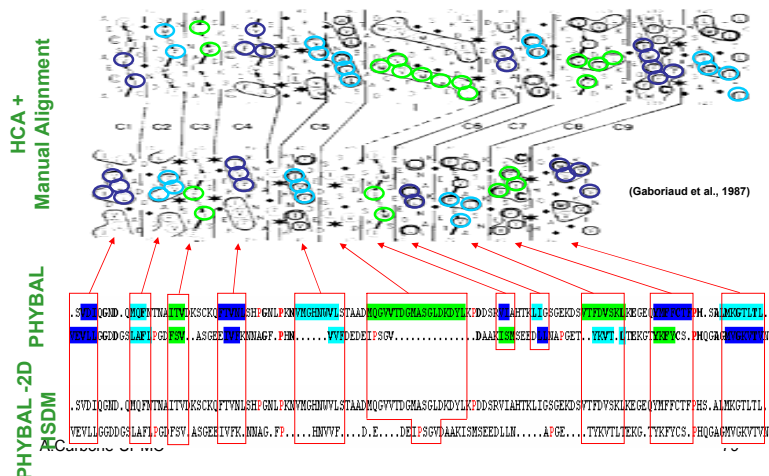
improve pairwise alignment
of distantly related proteins

A.Carbone-UPMC

78

Comparison with Hydrophobic Cluster Analysis method

Plastocyanin – Azurin (β, 13% Id)



Comparison with HMMSUM

Method		≤ 12	12-15	15-17	17-20	20-22	22-25
PHYBAL hot-spot	PCAP	44.9	39.4	40.3	52.8	57.9	75.0
	CAP	257	1037	1231	2189	1054	3737
	glo	17	18	78	93	167	74
	loc	2	2	22	41	80	42
HMMSUM	PCAP	38.8	35.1	38.0	54.3	59.2	76.9
	CAP	200	927	1278	2062	1138	3836
	glo	0	4	5	43	24	43
	loc	12	8	5	16	14	20

HMMSUM predicts secondary structures and aligns by using 281 structural context-based substitution matrices

80

(Huang and Bistoff, 2006)

Back to the motivations

- Detection of homologous proteins
- Detection of protein binding sites
- Multiple alignment

Weak Homology : some difficult cases

The length-dependency problem

```

1fel . AQEPVKGQVSTKPGSCP LILIRCAMLNPPNRCLKD.TDCPGIKKCCGEGSCGHACFVPQ...
1a0aA MKRESHKHAEQARRNRLAVLSELASLIPAEWQONVSAAPSKATTVEAACHRYIRHLQONGST
    
```



Small blocks of identities

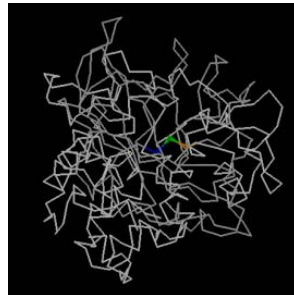
```

1fel AQEPVKGQVSTKPGSCP LILIRCAMLNPPN...RCLKDIDCFGIKKCCGEGSCGHACFVPQ
1udkA .....NEKSGSCPDLSH...PIFPLGICKTLCNSDSGGPNVQKCCCKNGCGFMTCTTPVP
    
```

Non homologous with some small blocks of identities

```

1cnsA SSSIVSRAQFDRHLHFDGACQAKFTYDAFAAAASSGTGSADVQKREVAAFLAQTSHE.TTGGWATAPDGAAVGYCRKQERGSSDYCT
1j3A ...QPTAPKDESGFNNDGTTQ...GFSV...NPDSPETAINENANN...ASSNSNSKGSNDLSEGGWAVVRISADING...GSINITYGD..I
PSAQWPCAPKFRWGRGFTCLSSNYNTGPAGRAGVDLLAMEDLVATDATVSFKTAWFWHTAQPPKFSSHALVCGCSSPSGADRAAGRVPGWWITLSIS
KLTDVIAP.....TPAVS...AVLPQSSTHG...UGH.....TRAIR...WTNNFVAQTDCTTALLTESTNDSPNAVATDAADSVYWNI
VSGECPHGDSEVADHIGFYKRYCDILGVGYGNLDCVSORPFA
...LFVSNSDNESLDWIATK
    
```



Similar score to weakly homologous proteins

Homologous with no block of identities

```

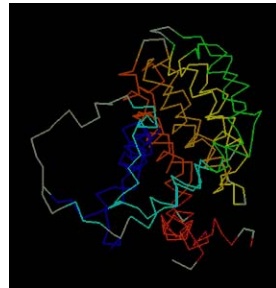
1bwwA TEDWQSPSSLSASVGRVITTCQASDDLIKYLWWQKPGRAPKLLLEASNLQAGVPSHFGSGSDTVYFTISHQPEDLAVYCOQYSHF
1cdA .....RDSGTWGADIGDLANFNQNDDIDEVW...ERGSTVAFYRKRGPK...SGAFEILANG...DLIAWWTRDSGTEWYSTNGT
VYFCOQTRLSTT
RILDKALDLWLE
    
```



Non-homologous : but partial homology

```
Ik3aA ..... N E S T L . . P G E V L A I R G I F H A G G I N . E E E T L T H P . S P I K L Y I T S L I R D K E S . . . . . F L I M L A N Y R E H S T T C I I D O M O S M K W S G D G
IghA  MAHAGRTGYDNEF E L A V I E Y L S O R G Y E D A G I D V E E N R T E A P E G T E S E T V P L T L R Q A G D D S G Y P I G L I S S O L L I P P T T A R G P F A T V V E E L F D G
      . . . . .
N I D W G R L A I L T F G S . F V A O K L S M E . P E L S D F A L A V P A I V E X M G P O T T R A R G G V S G L K A Y C T O W I T D D D L E H H H H H H
. V N W G R I V A F F F G G V M C V E S V N R E M S P L D N I A L V R T Y L N R H L E L V Q D N G G L A F V E L G P S M R . . . . .
```

Globin-like domain in common



A.Carbone-UPMC

Références bibliographiques

D.S.Hirschberg, Algorithms for the longest common subsequence problem, *J. ACM*, 24:664-675, 1977.

D.Gusfield, *Algorithms on strings, Trees and Sequences*, Cambridge University Press, 1997.

J.Baussand, C.Deremble, A.Carbone, Periodic distributions of hydrophobic amino acids allows the definition of fundamental building blocks to align distantly related proteins, *Proteins: Structures, Functions and Bioinformatics*, 2006.

A.Carbone-UPMC

86

Autres références

P.Clote et R.Backofen, *Computational Molecular Biology*, Wiley Series, Mathematical and Computational Biology, Princeton University Press, 2001.

R.Durbin, S.Eddy, A.Krogh, G.Mitchison, *Biological Sequence Analysis, probabilistic models of proteins and nucleic acids*, Cambridge University Press, 2000.

P.Pevzner, *Computational Molecular Biology, an algorithmic approach*, MIT Press, 2000.

D.Graur et Wen-Hsiung Li, *Fundamentals of Molecular Evolution*, Sinauer, second edition, 1999.

I.Kirching, P.Forey, Ch. Humphries, D.Williams, *Cladistics, the theory and practice of parsimony analysis*, Oxford University Press, second edition, 2000.

C.Branden et J. Tooze, *Introduction to Protein Structure*, Garland Publishing, second edition, 1999.

A.Carbone-UPMC

87