

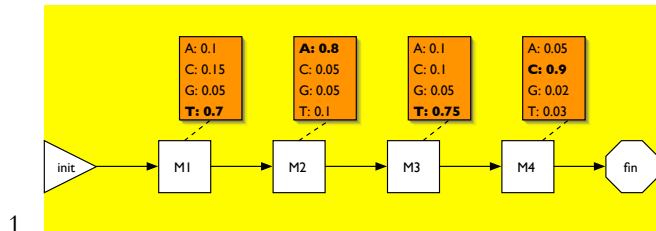
5 Applications de HMMs

5.1 Motif de séquence conservé

Un alignement multiple de protéines peut révéler des motifs conservés :

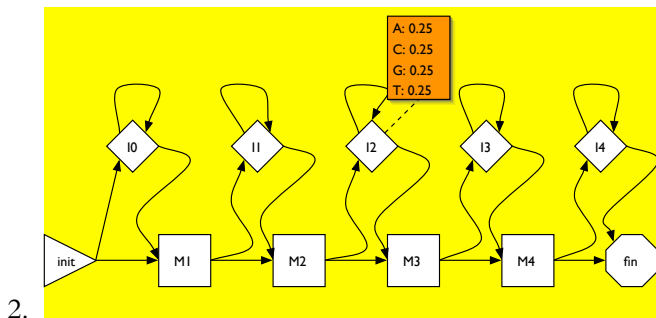
Q9DAK4 /4-132	QLQG.TWKSVCDFN.FENYMKELGVGRASRK.LGCLAK.....PTVT
Q8QHA8 /4-132	HFVG.TWKLLSSENF.EDYMKELGVGFATRK.MAGVAK.....PNLT
FABA_RAT /3-131	AFVG.TWKLVSSSENF.DDYMKEVGVGFATRK.VAGMAK.....PNLI
MYP2_HUMAN /3-131	KFLG.TWKLVSSSENF.DDYMKALGVGLATRK.LGNLAK.....PTVI
TLBP_MOUSE /4-132	PFLG.TWKLVSSSENF.ENYVRELGVCEPRK.VACLIK.....PSVS
Q57663 /4-133	KFVG.TWKMISSDNF.DDYMKAIQGVGFATRQ.VGNRTK.....PNLV
FABE_RAT /6-134	DLEG.KWRLVESHGF.EDYMKELGVGLALRK.MGAMAK.....PDCI
FABE_BOVIN /6-134	QLVG.RWRLVESKGF.DEYMKEVGVGMALRK.VGAMAK.....PDCI

HMM de profile. On construit un modèle HMM dont les états correspondent aux colonnes de l'alignement :



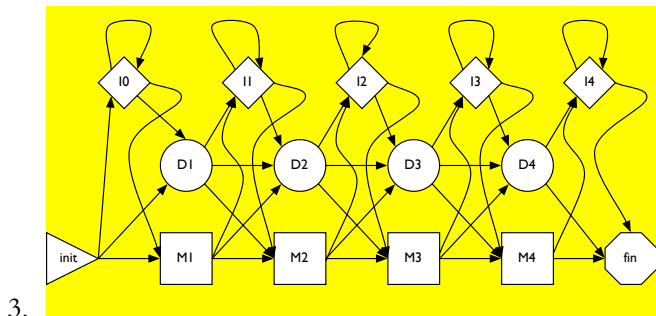
1.

substitutions : colonnes de résidus alignés



2.

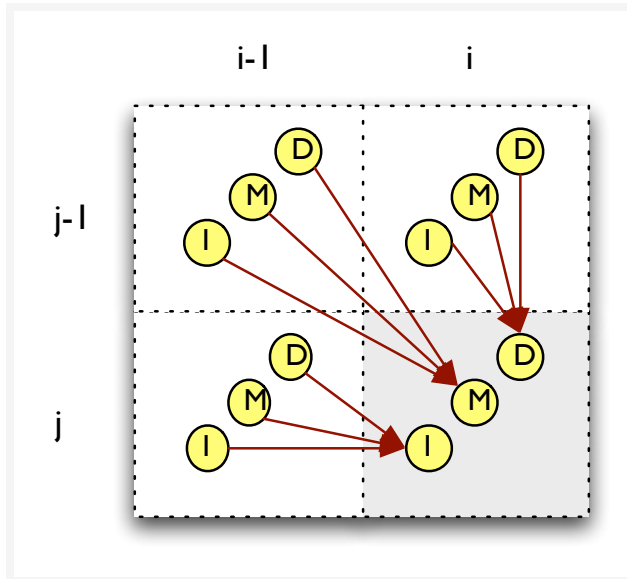
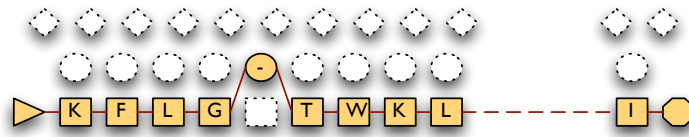
insertions : résidus entre colonnes conservées



3.

suppressions : états sans émission !

Alignement à un profil. On calcule le chemin Viterbi pour aligner une nouvelle séquence au profil



Programmation dynamique :

- ★ états de match : $M_j: j = 1, \dots, n$
- ★ états d'insertion : $I_j: j = 0, \dots, n$
- ★ états de suppression : $stD_j: j = 1, \dots, n$
- ★ sous-problèmes pour PD : $V^{(M)}(i, j), V^{(I)}(i, j), V^{(D)}(i, j)$ pour émission de $s_1 \dots s_i$ en finissant dans un état M_j, D_j, I_j

Pondération par profil. Hypothèse null : $S[1..\ell]$ est une séquence de caractères iid $\mathbb{P}\{S[i] = a\} = r_a$ (souvent les mêmes probs d'émission aux états d'insertion). Hypothèse alternatif : $S[1..\ell]$ est générée par notre HMM de profil \mathcal{M} . LODS (logarithme des chances) : comparer $p_0 = \mathbb{P}\{S[1..\ell] \mid H_0\}$ et $p_{\mathcal{M}} = \mathbb{P}\{S[1..\ell] \mid H_1\}$:

$$\text{LODS} = \log \frac{p_{\mathcal{M}}}{p_0}$$

Calcul de LODS : par Viterbi, remplacer probabilité d'émission $p_q(a)$ par $p_q(a)/r_a + \text{logarithme}$

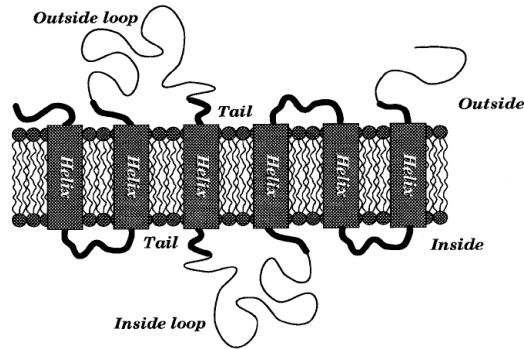
Bases de données de profils. PFAM et PROSITE utilisent des HMMs de profil pour des familles de protéines. Exemple : <http://prosite.expasy.org/PS50071>. On trouve la pondération de substitutions par colonne, ainsi que la pondération de transitions entre états (en LODS).

```
ID HOMEBOX_2; MATRIX.
...
      A  B  C  D  E  ...
...
/M: SY='E'; M= -5,  2,-25,  3, 11, ...
/M: SY='Q'; M= -3, -4,-25, -4, 12, ...
...
/I:      I=-8; MI=-8; IM=-8; DM=-15; MD=-15;
```

5.2 Architectures spécialisées

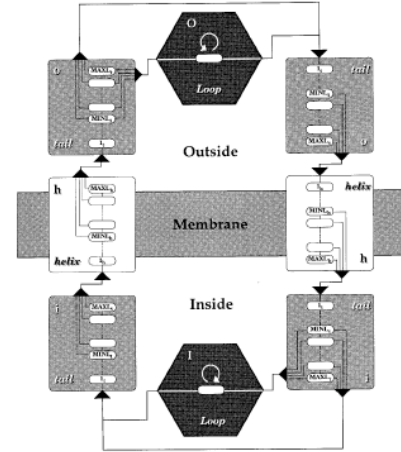
L'architecture de l'HMM (transitions permises entre états) peut capturer des contraintes spécifiques.

Protéines transmembranaires. On définit des états pour l'intérieur/extérieur de membrane.



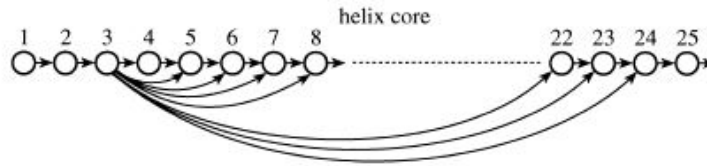
Amino acid seq: MGDVCDTEFGILVA...SVALRPRKHGRWIV...FWVDNGTEQ...PEHMTKLHMM...
State seq: ooooooooohhhh...hhhhiiiiihhh...hhhooooO...OOoooohh...

(Tusnady & Simon, J.Mol.Biol. 283:489)



(Tusnady & Simon, J.Mol.Biol. 283:489)

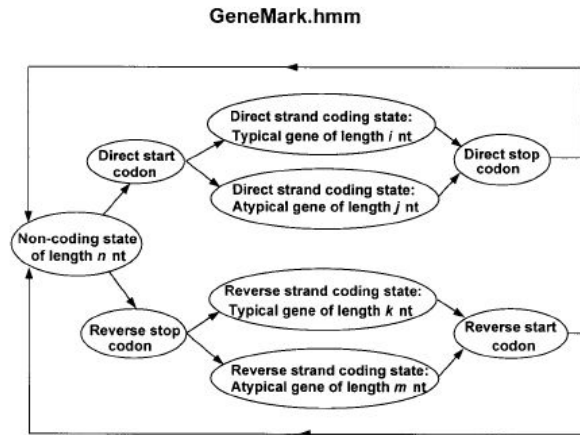
Modèles de durée. durée de séjour dans un état : p.e., longueur du segment restreinte pour hélices dans le membrane (min et max) :



5.3 Gènes

Les états de l'HMM modèlisent des composantes d'un gène.

GeneMark. Gènes compacts chez procaryotes. . .



GENSCAN. Eucaryotes : modèles plus compliqués.

