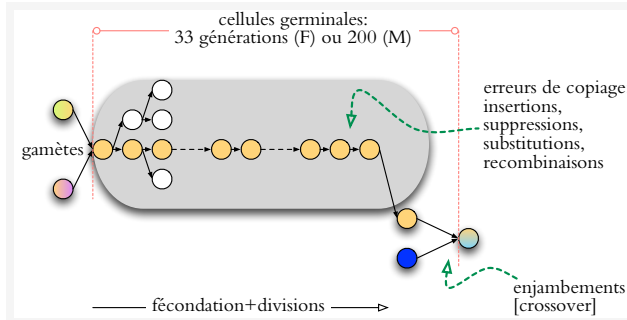
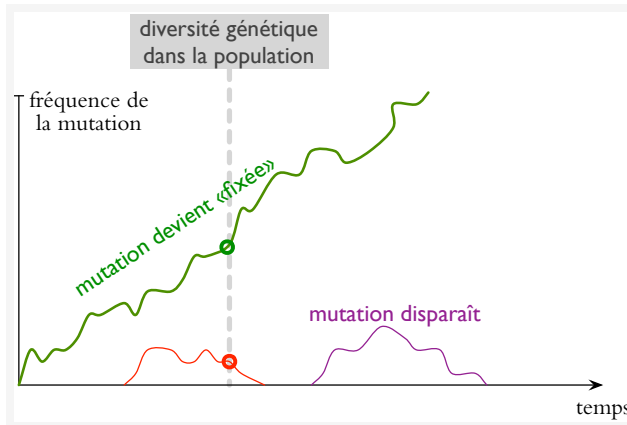


6 Modèles temporels

6.1 Descendance avec modification



Les mutations lors d'héritage du matériel génétique entre générations introduisent de nouveaux allèles dans une population. Un allèle du parent (p.e., mutation nouvelle dans la lignée de cellules germinales) peut être hérité par l'enfant.



Ultimement, l'allèle devient fixé dans la population entière, ou disparaît complètement. La probabilité de fixation dépend (a) de la taille de la population, et (b) l'effet de l'allèle (avantage ou désavantage possible). Les mutations fixées s'accumulent → divergence de séquences d'origine commune. La divergence croît avec le nombre de générations depuis l'ancêtre.

Substitutions multiples. Évolution de caractères homologues est **sans mémoire** : états successifs forment une chaîne de Markov. Probabilités de **substitution** selon la matrice $\mathbf{M}_{X \rightarrow Y}$ où

$$\mathbf{M}_{X \rightarrow Y}[a, b] = \mathbb{P}\{Y = b \mid X = a\}$$

avec valeurs $a, b \in \Sigma = \{1, \dots, r\}$. (Exemple : $\Sigma = \{1, 2, 3, 4\}$ encodant ADN)

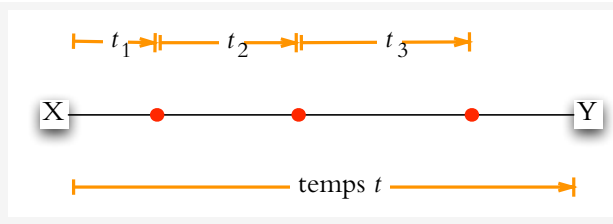
$$\begin{array}{c} X \\ \downarrow \\ Y \\ \downarrow \\ Z \end{array}$$

$$\begin{aligned} \mathbb{P}\{Z = z \mid X = x\} &= \sum_{y=1}^r \mathbb{P}\{Z = z, Y = y \mid X = x\} = \sum_{y=1}^r \mathbb{P}\{Z = z \mid X = x, Y = y\} \mathbb{P}\{Y = y \mid X = x\} \\ &= \sum_{y=1}^r \mathbb{P}\{Z = z \mid Y = y\} \mathbb{P}\{Y = y \mid X = x\} \quad (\text{propriété de Markov}) \end{aligned}$$

Donc, on peut écrire que $\mathbf{M}_{X \rightarrow Z} = \mathbf{M}_{X \rightarrow Y} \cdot \mathbf{M}_{Y \rightarrow Z}$

6.2 Temps continu.

Comment peut-on ajouter une notion de temps ? Quelle est la matrice \mathbf{M} qui correspond à un temps de divergence connu entre les séquences ?



À chaque mutation, le changement d'état est déterminé par la matrice \mathbf{M} . Si k événements en temps t , alors $\mathbf{M}_{X \rightarrow Y} = \mathbf{M}^k$. Donc la matrice pour temps t devrait être $\mathbf{M}^{N(t)}$ où $N(t)$ est le nombre de mutations pendant temps t . Problème : $N(t)$ est aléatoire. . .

Processus de Poisson. Supposons que $\{N(t) : t \in [0, \infty)\}$ est un processus qui compte l'occurrence d'événements de quelque sort. Supposons que le processus satisfait les critères suivants

☞ [«orderly»] La probabilité qu'il y ait plus d'une occurrence dans un petit intervalle de temps est négligeable

$$\lim_{\delta \rightarrow 0} \mathbb{P}\{N(t + \delta) > 1 \mid N(t + \delta) \geq 1\} = 0$$

☞ [«memoryless»] Le nombre d'occurrences pendant un intervalle ne dépend pas des événements précédents : $N(t + s) - N(t)$ est indépendant de $N(t)$

Si le processus est *orderly* et *memoryless*, alors c'est un **processus de Poisson** avec les propriétés suivantes.

* nombre d'arrivées pendant temps s est une v.a. Poisson :

$$\forall t, s : \mathbb{P}\{N(t + s) - N(t) = k\} = e^{-\lambda s} \frac{(\lambda s)^k}{k!}$$

Le paramètre λ est l'intensité ou le **taux** du processus.

* temps d'attente est une v.a. exponentielle : $\forall t, s : \mathbb{P}\{N(t + s) = N(t)\} = e^{-\lambda t}$

Processus de Markov. En retournant à ce qui se passe pendant temps t :

$$\sum_{k=0}^{\infty} e^{-\lambda t} \frac{(\lambda t)^k}{k!} \mathbf{M}^k = \exp(\lambda t(\mathbf{M} - \mathbf{I})) = \exp(\lambda t \mathbf{Q})$$

où \mathbf{I} est la matrice d'unité. $\mathbf{Q} = \mathbf{M} - \mathbf{I}$ est la *matrice instantanée de taux de substitutions*. Calcul de $\exp(\lambda \mathbf{Q}t)$: décomposition de la matrice $\mathbf{Q} = \mathbf{U} \mathbf{\Lambda} \mathbf{V}$, et faire $e^{\lambda \mathbf{Q}t} = \mathbf{U} e^{\lambda t \mathbf{\Lambda}} \mathbf{V}$.

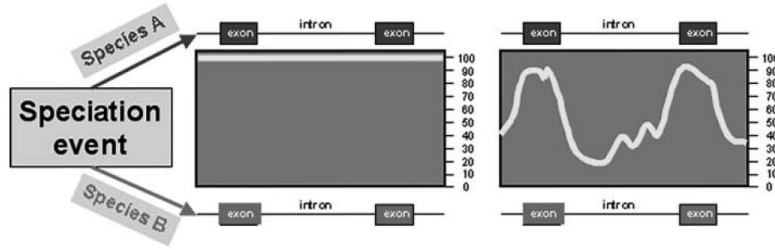
$$\begin{pmatrix} -\eta & \eta \\ \nu & -\nu \end{pmatrix} \begin{pmatrix} \cdot & \frac{\mu}{4} & \frac{\mu}{4} & \frac{\mu}{4} \\ \frac{\mu}{4} & \cdot & \frac{\mu}{4} & \frac{\mu}{4} \\ \frac{\mu}{4} & \frac{\mu}{4} & \cdot & \frac{\mu}{4} \\ \frac{\mu}{4} & \frac{\mu}{4} & \frac{\mu}{4} & \cdot \end{pmatrix} \begin{pmatrix} \cdot & \pi_{\text{C}\alpha} & \pi_{\text{G}\beta} & \pi_{\text{T}\alpha} \\ \pi_{\text{A}\alpha} & \cdot & \pi_{\text{C}\alpha} & \pi_{\text{T}\beta} \\ \pi_{\text{A}\beta} & \pi_{\text{C}\alpha} & \cdot & \pi_{\text{T}\alpha} \\ \pi_{\text{A}\alpha} & \pi_{\text{C}\beta} & \pi_{\text{G}\alpha} & \cdot \end{pmatrix} \quad (\text{somme est } 0 \text{ dans toute rangée})$$

perte-gain (binaire) Jukes-Cantor (DNA) Hasegawa-Kishino-Yano (DNA)

Le modèle de Yang et Nielsen (1998) pour codons utilise un paramètre ω qui mesure sélection, et un paramètre κ pour rapport de transitions/transversions :

$$\mathbf{Q}[a, b] = \begin{cases} \mu\pi_b & 1 \text{ transversion synonyme} \\ \mu\kappa\pi_b & 1 \text{ transition synonyme} \\ \mu\omega\pi_b & 1 \text{ transversion non-synonyme} \\ \mu\kappa\omega\pi_b & 1 \text{ transition non-synonyme} \end{cases}$$

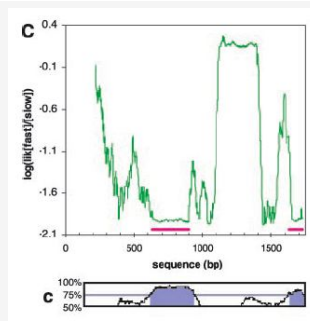
6.3 Conservation de séquence



Miller & al. *Annu Rev Genomics Hum Genet* 5 :15 (2004)

Principe de génomique comparative : éléments fonctionnels sont plus conservés (sélection négative) que les éléments non-fonctionnels (évolution neutre)

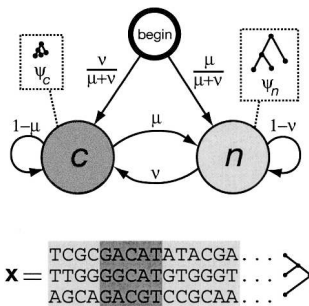
Phylogenetic shadowing. Méthode : modèle de substitutions dans les cas de sélection négative et neutre



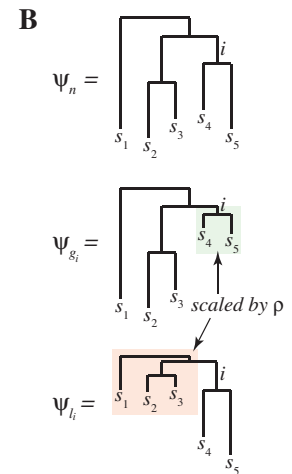
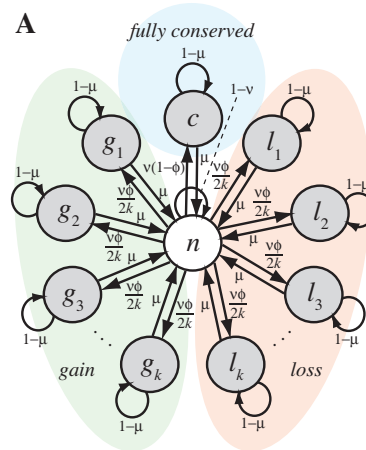
Comparaison de séquences entre des espèces proches : évolution rapide/lente $e^{\mu Q t}$ où $\mu < 1$ pour des régions sous sélection purificatrice.
Exemple : fenêtre glissante de 50 pb sur alignement multiple de primates.
Axis Y : score LODS (en haut) et % d'identité (en bas)

Boffelli & al. *Science* 299 :1391 (2003)

PhyloHMM. Modèle HMM pour segmentation : émissions = colonnes de l'alignement multiple, avec probabilités $e^{Q t}$ (neutre) ou $e^{\rho Q t}$ (sélection négative avec $\rho < 0$)



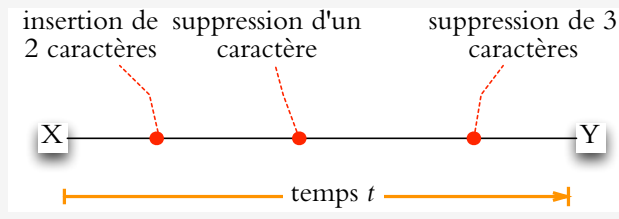
deux états : conservé ou non
Siepel & al. *Genome Res* 15 :1034 (2005)



DLESS : états correspondent à conservation dans un clade phylogénétique, ou à la perte de fonction
Siepel & al. *RECOMB* 2006

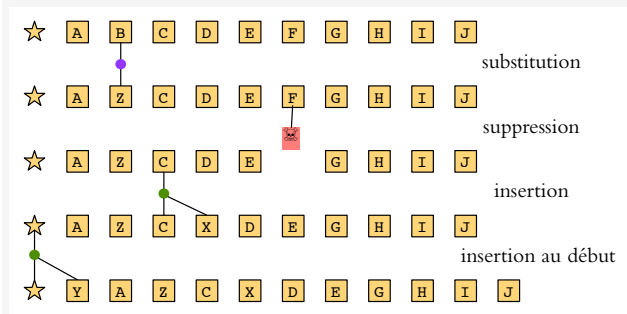
6.4 Modèle temporaire pour suppressions et insertions

On a vu le modèle de Markov pour substitutions — est-ce qu'on peut modéliser les indels par un processus stochastique ?

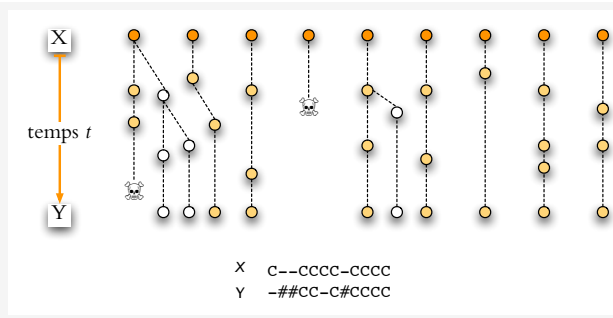


Idee : événements forment un processus (p.e., Poisson de taux θ), il faut juste spécifier ce qui se passe lors d'un événement. Un modèle assez général : insertions de longueur ℓ avec probabilité $p_{Ins}(\ell)$ et suppressions de longueur ℓ avec probabilité $p_{Del}(\ell)$ — très difficile (ou impossible) de travailler avec...

Thorne-Kishino-Felsenstein [1991] Un modèle simple : insertions et suppressions de longueur 1.



Un processus indel (en parallèle avec les processus de substitution) à chaque caractère + au début. Événement d'insertion : choix d'un caractère selon la distribution stationnaire. Taux d'insertion λ à chaque caractère ainsi qu'au début ; taux de suppression μ à chaque caractère.



But : étant donné une séquence ancestrale et une séquence descendante, établir la homologie au niveau des caractères (représentée par un alignement). Les caractères s'évaluent indépendamment : on peut considérer le **sort** de chaque caractère ancestral séparément.

Sort	Description	Bloc d'alignement
$C \underbrace{\# \dots \#}_{n \text{ fois}}$	ancêtre survivant, avec $n \geq 0$ caractères insérés à côté	$\begin{matrix} C--- \\ C### \end{matrix}$
$\square \underbrace{\# \dots \#}_{n \text{ fois}}$	ancêtre mort, avec $n > 0$ caractères insérés à côté	$\begin{matrix} C--- \\ -### \end{matrix}$
\square	ancêtre mort, aucune insertion à côté	$\begin{matrix} C \\ - \end{matrix}$
$\star \underbrace{\# \dots \#}_{n \text{ fois}}$	insertion de $n \geq 0$ caractères au début	$\begin{matrix} --- \\ ### \end{matrix}$

Sort	probabilité $X \rightarrow Y$
$C \rightarrow C \underbrace{\# \dots \#}_{n-1}$	$H_t B_t^{n-1}$
$C \rightarrow \square \underbrace{\# \dots \#}_{n}$	$N_t B_t^{n-1}$
$C \rightarrow \square$	E_t
$\star \rightarrow \star \underbrace{\# \dots \#}_{n}$	$I_t B_t^n$

$$\beta_t = \frac{1 - e^{-(\mu - \lambda)t}}{\mu - \lambda e^{-(\mu - \lambda)t}}$$

$$B_t = \lambda \beta_t \qquad E_t = \mu \beta_t$$

$$I_t = 1 - \lambda \beta_t \qquad H_t = e^{-\mu t} (1 - \lambda \beta_t)$$

$$N_t = (1 - e^{-\mu t} - \mu \beta_t)(1 - \lambda \beta_t);$$

Maintenant on peut

- ★ trouver le meilleur alignement
- ★ maximiser la vraisemblance pour deux séquences homologues et ainsi trouver les valeurs de λ, μ, t
- ★ concevoir des tests de homologie basés sur la valeur de la vraisemblance
- ★ évaluer la fiabilité de l'alignement optimal

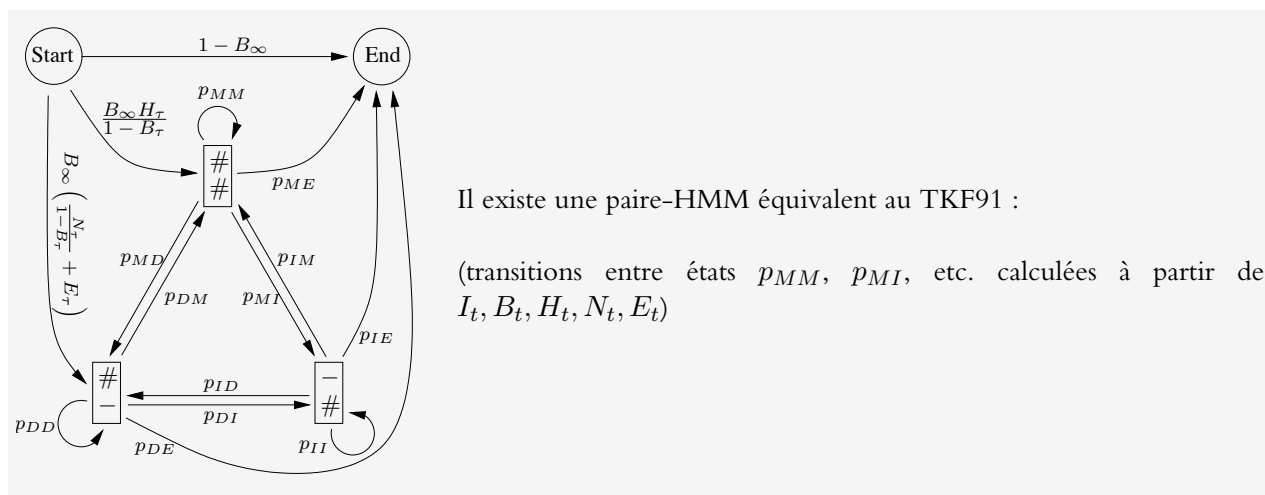
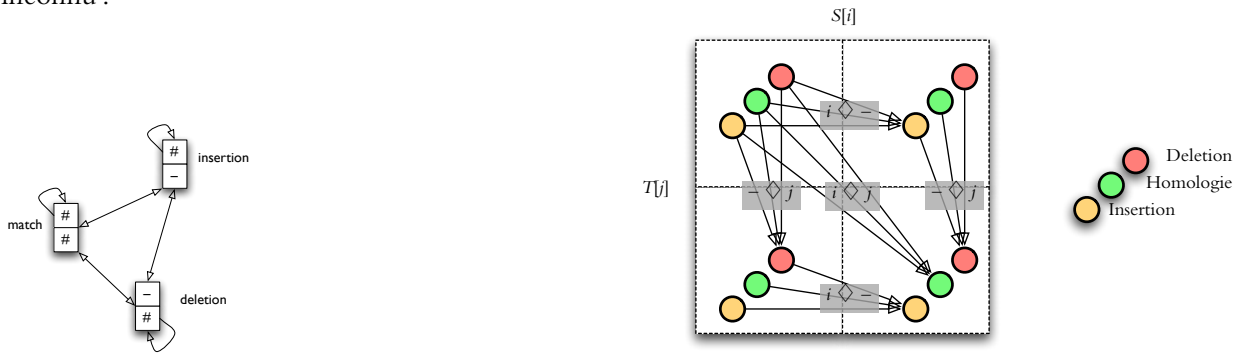
Exemple : calcul de la vraisemblance — utiliser la programmation dynamique (en se servant des queues géométriques $\propto B^n$). $L(i, j)$ probabilité que $X[1..i]$ est devenu $Y[1..j]$; $Z(i, j)$ variables auxiliaires. Récurrences :

$$Z(i, j) = p_{i,j} \cdot H_t \cdot Z(i - 1, j - 1) + p_j \cdot N_t \cdot Z(i, j - 1) + p_j \cdot B_t \cdot Z(i, j - 1)$$

$$L(i, j) = Z(i, j) + E_t \cdot L(i - 1, j),$$

où $p_{i,j} = \mathbf{M}[X[i], Y[j]]$ est la probabilité $X[i] \rightarrow Y[j]$ dans la matrice de substitution, et $p_j = \pi_{Y[j]}$ est la probabilité stationnaire de $Y[j]$. Initialisation : $Z(0, 0) = L(0, 0) = I_t$.

Paire-HMM. Modèle de Markov caché où les états correspondent à des insertions, suppressions et substitutions/identités dans l'alignement. La machine *émet* un alignement, mais on observe seulement les séquences sans trous. Problème d'inférence : on connaît les séquences — qu'est-ce qu'on peut dire sur l'alignement inconnu ?

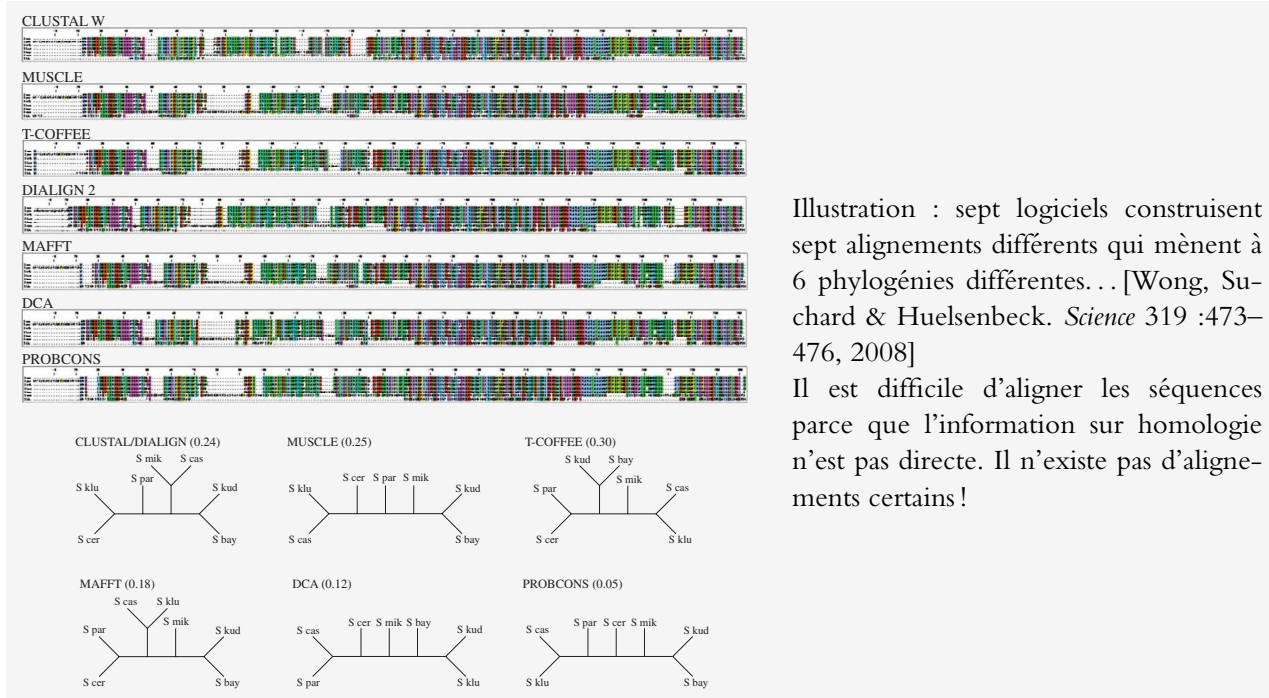


Il existe une paire-HMM équivalent au TKF91 :

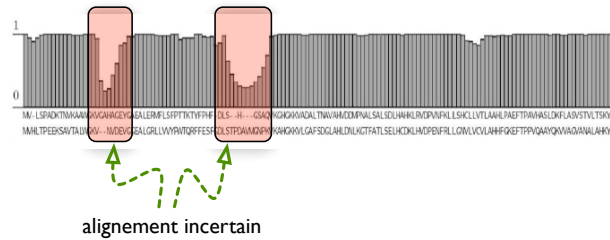
(transitions entre états p_{MM} , p_{MI} , etc. calculées à partir de I_t, B_t, H_t, N_t, E_t)

6.5 Alignement ?

alignement = hypothèse de homologies



Fiabilité. On a la probabilité postérieure de homologie ($i \diamond j$), et celle de manque de homologie ($i \diamond -$, $-\diamond j$) pour juger les colonnes de l'alignement. Une partie bien alignée correspond à une séquence de colonnes avec de grandes probabilités postérieures.



Maximal Expected Accuracy. l'alignement qui maximise le nombre de paires alignés correctement — utiliser $\mathbb{P}\{i \diamond j\}$ comme le score du match de la colonne $\frac{X[i]}{Y[j]}$, pas de pénalité pour trous.

AMAP. Problème avec MEA — qu'est-ce qui se passe avec deux séquences non-relées ? Maximiser plutôt le nombre (en espérance) d'accords entre le vrai alignement (inconnu) et l'alignement qu'on calcule :

$$\text{désaccords} = n + m - 2M - D - I,$$

où M, D, I sont les nombres de matches, suppressions et insertions en commun entre deux alignements. En espérance, on a $\mathbb{E}M = \sum_{i,j} \mathbb{P}\{i \diamond j\}$, $\mathbb{E}D = \sum_i \mathbb{P}\{i \diamond -\}$, $\mathbb{E}I = \sum_j \mathbb{P}\{-\diamond j\}$ (où \mathbb{P} dénote probabilité postérieure, conditionnée sur les deux séquences).