

PRÉDICTION DE GÈNES

GÈNES

Question principale : quel genre de gènes ?

Procaryotes [pas d'introns, opérons] et Eucaryotes [introns, signalisation compliquée]

Gènes traduits en protéines

Gènes non-traduits (gènes ARN)

Gènes non-transcrits (signaux de réplication, recombinaison, ségrégation [meiosis], ...)

RAPPEL : DE L'ADN À PROTÉINE

1. transcription : copie du brin informatif à ARN messenger
(ARN : utilise Uracile au lieu de Thymine)
2. traduction : ARNm à protéine (par le ribosome) : acides aminés fournis par ARN de transfert

TERMINOLOGIE

similarité : notion algorithmique de relation entre séquences

homologue : relié par un ancêtre commun

→ **orthologue** : relié par événement de spéciation

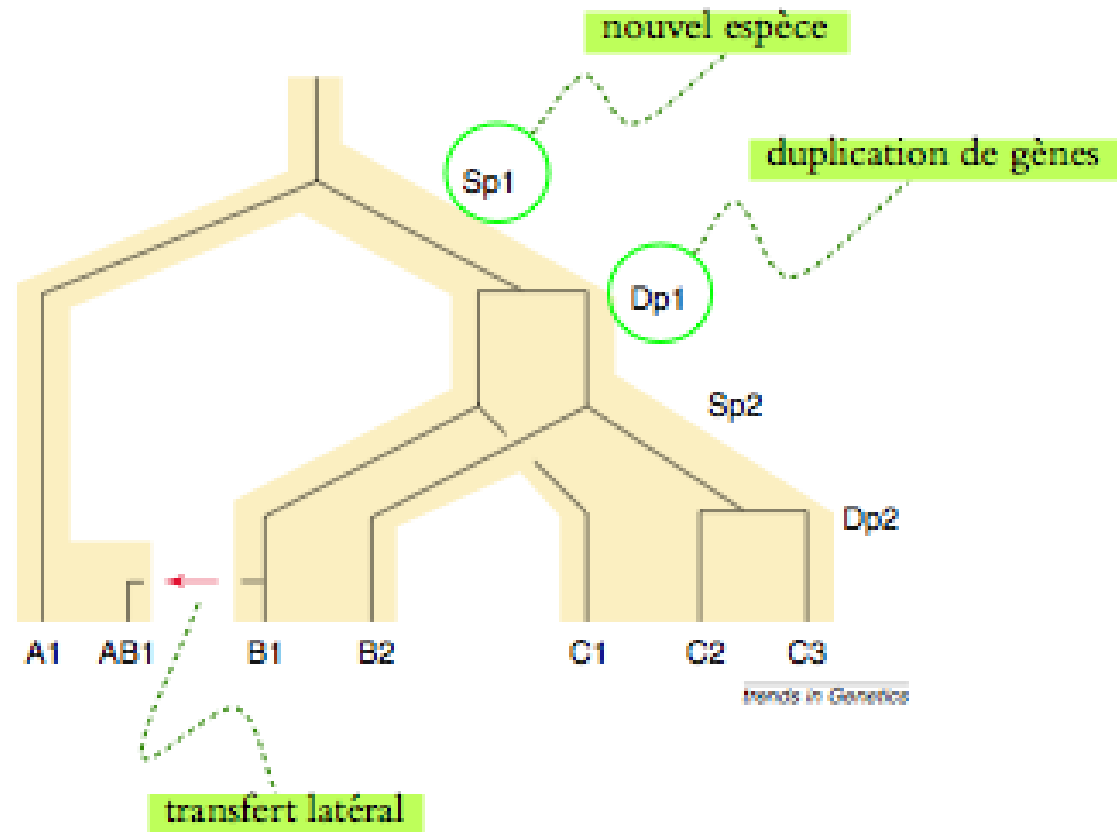
→ **paralogue** : relié par événement de duplication

→ **xenologue** : acquis par un autre mécanisme (transfert latéral)

similarité n'implique pas toujours la homologie : évolution convergente

homologie n'implique pas toujours la similarité non plus. . .

GÈNES HOMOLOGUES



orthologues : B1–A1, B1–C1

paralogues : B1–B2 (*in-paralog*), B1–C2 (*out-paralog*)

xenologues : A1–AB1 co-orthologues : {C1, C2, C3}–{B1, B2}

Fitch *Trends in Genetics* 16 :227 (2000)

GÉNOMIQUE COMPARATIVE

Alignement de deux séquences :

- évidence de homologie
- conservation indique la fonctionnalité
- étudier les mécanismes de mutation
- étudier les forces d'évolution

p.e. comparer le taux de mutations synonymes (entre codons encodant le même acide aminé) et celui de mutations non-synonymes

évolution neutre : aucune différence

évolution/sélection purificatrice/négative : synonyme plus fréquent

sélection positive [évolution Darwinien] : non-synonyme plus fréquent

GÈNES TRADUITS — PROCARYOTES

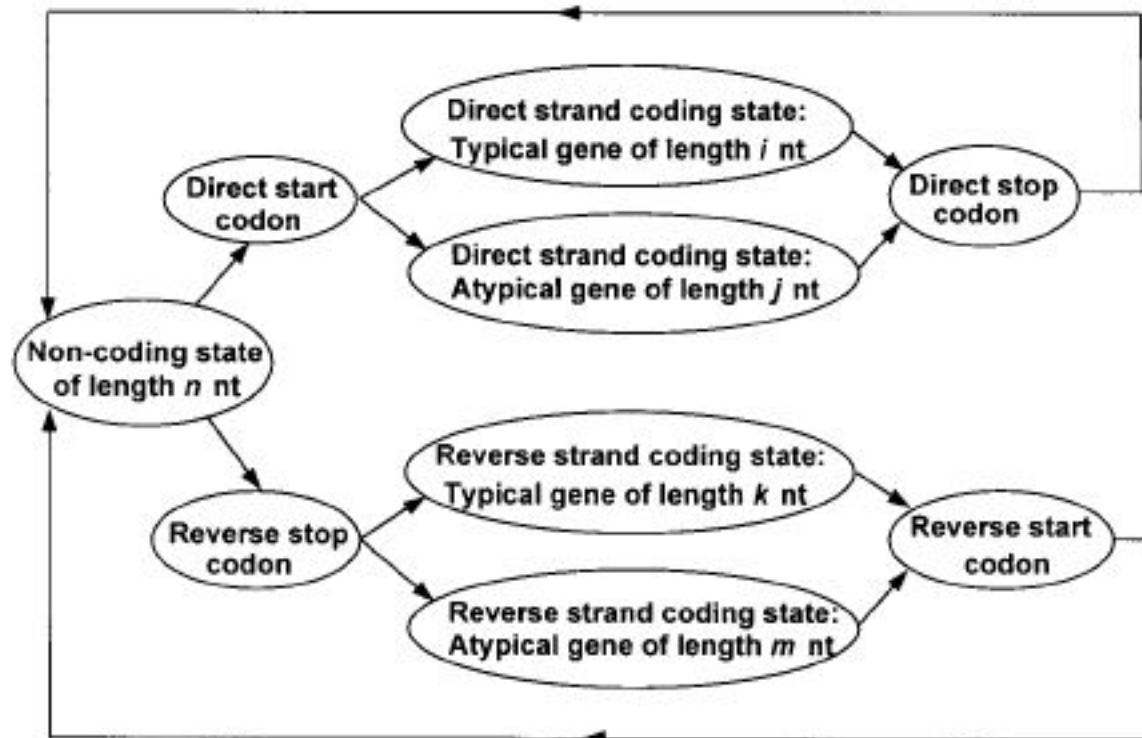
ORF : Open Reading Frame

ATG...Ter où $\text{Ter} \in \{\text{TAA}, \text{TAG}, \text{TGA}\}$

GÈNES TRADUITS — PROCARYOTES

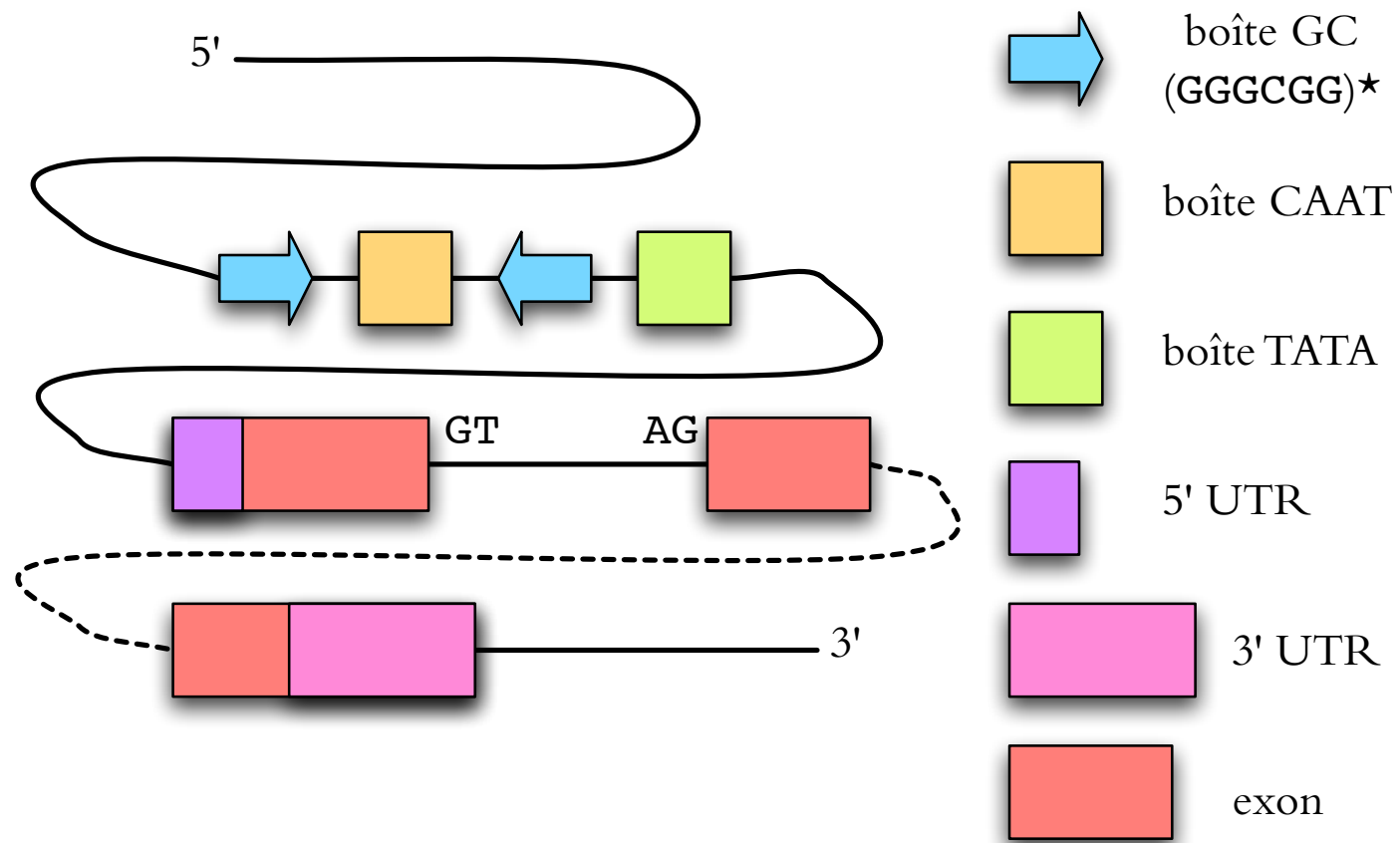
Procaryotes (pas d'exons !) — GeneMark

GeneMark.hmm



Lukashin & Borodovsky *Nucleic Acids Res* 26 : 1107 (1998)

GÈNES TRADUITS — EUCARYOTES

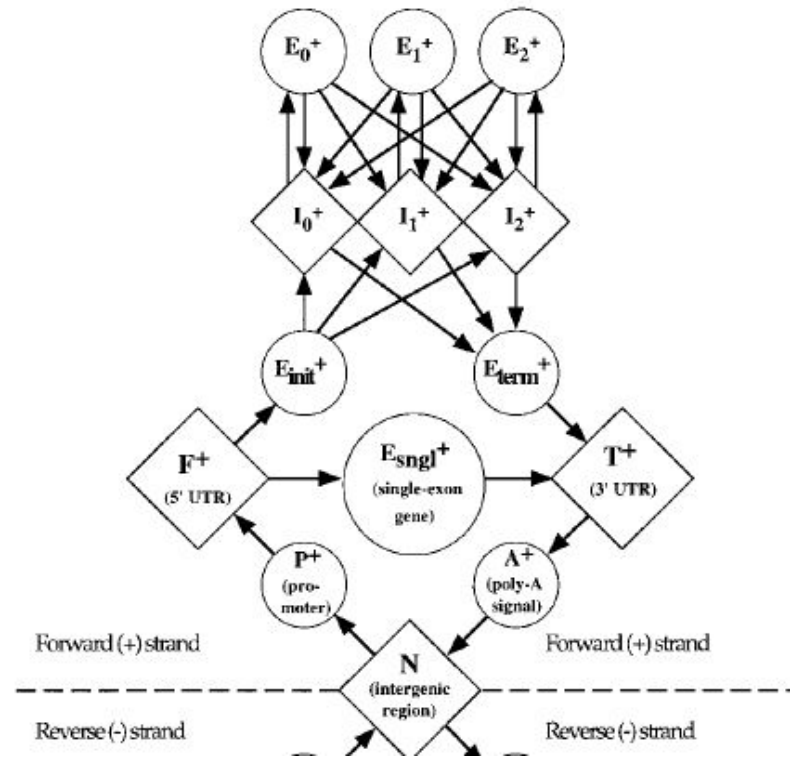


[processes : spliceosome, ribosome]

après Graur & Li *Fundamentals of Molecular Evolution* (2000)

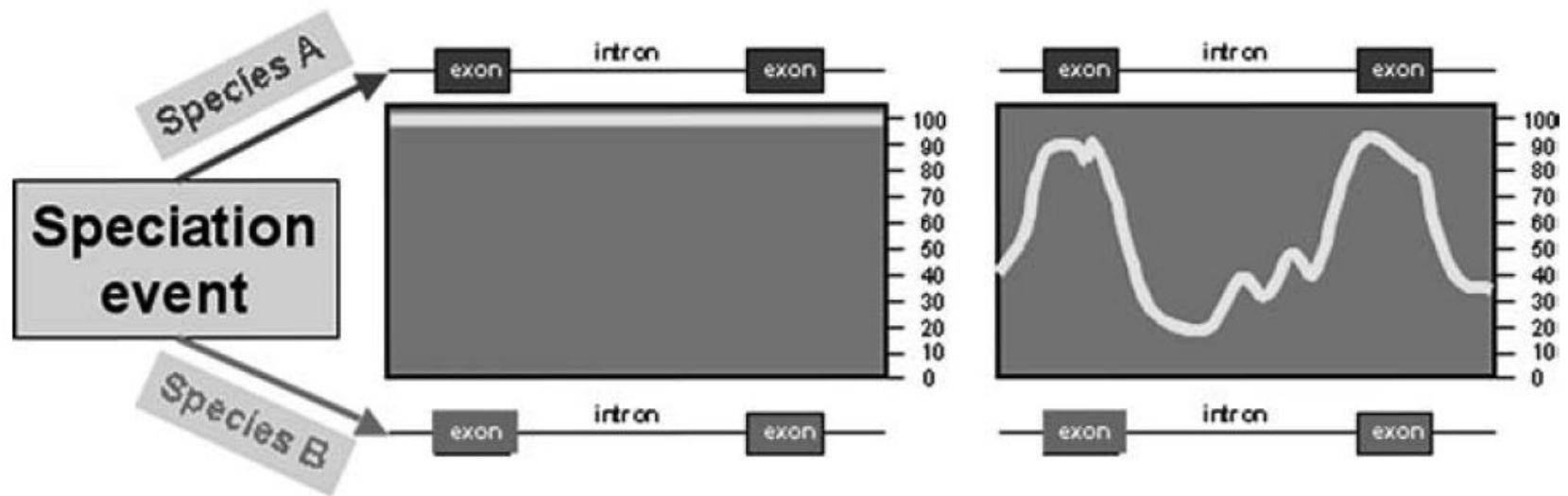
PRÉDICTION DE GÈNES (GENSCAN)

Eucaryotes — GENSCAN



Burge & Karlin *J Mol Biol* 268 : 78 (1997)

MÉTHODE COMPARATIVE POUR LA RECHERCHE DE GÈNES



Alignement de deux régions aide à l'identification d'exons : les exons sont plus préservés (sélection négative)

Principe de génomique comparative : éléments fonctionnels sont plus [sélection négative] ou moins [sélection positive] préservés que des éléments non-fonctionnels [évolution neutre]

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

PRÉDICTION DE GÈNES (SGP2)

séquences de référence
(p.e., génome d'un autre organisme)

alignements locaux

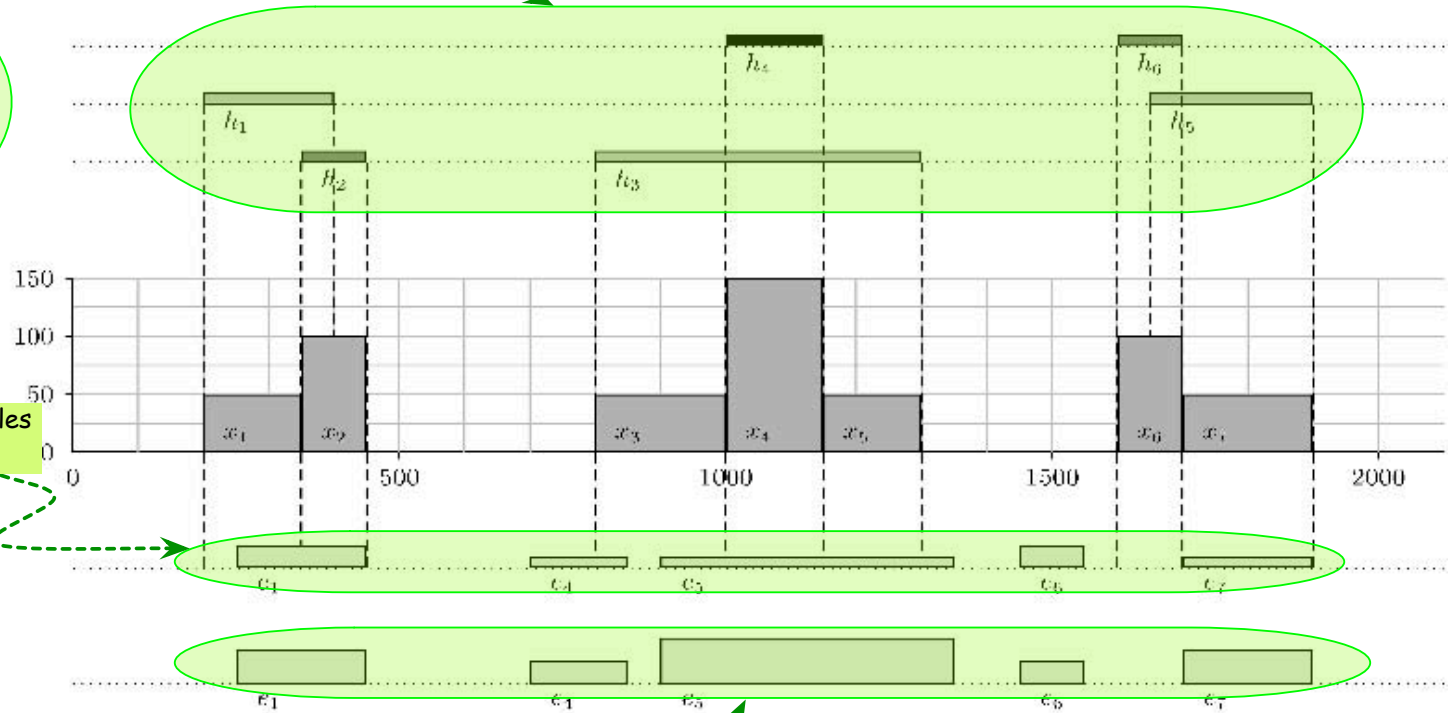
- A HSPs
- dbseq 1
 - dbseq 2
 - dbseq 3

B Max score of projected HSPs

prédictions initiales
(ab initio)

C EXONS

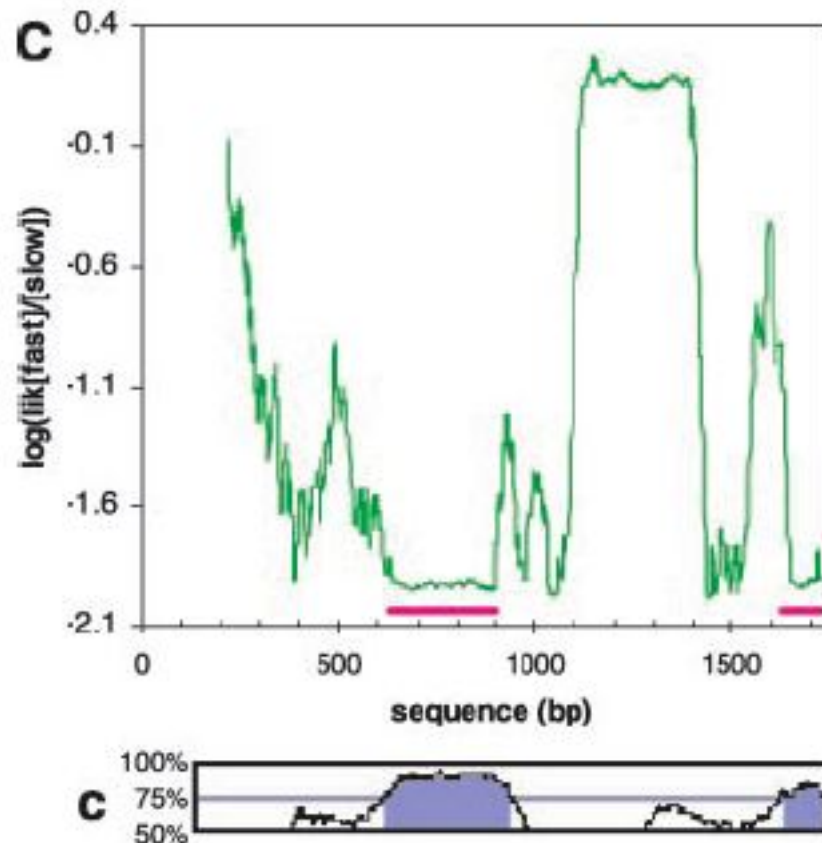
RESCORED EXONS



nouveau score combine
celui de la prédiction initiale et
celui des alignements

PHYLOGENETIC SHADOWING

Comparaison de séquences entre des espèces proches :
modèles d'évolution rapide/lente (HMM)

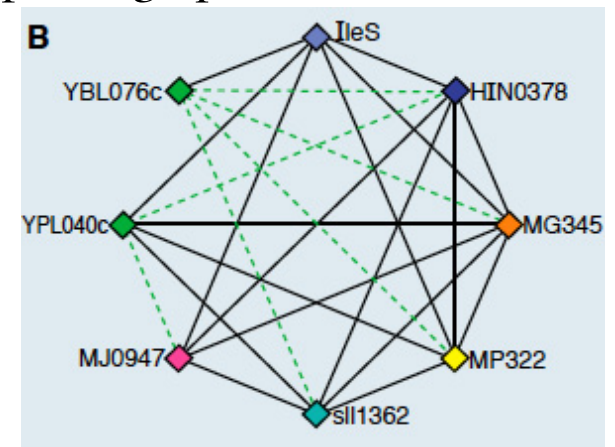
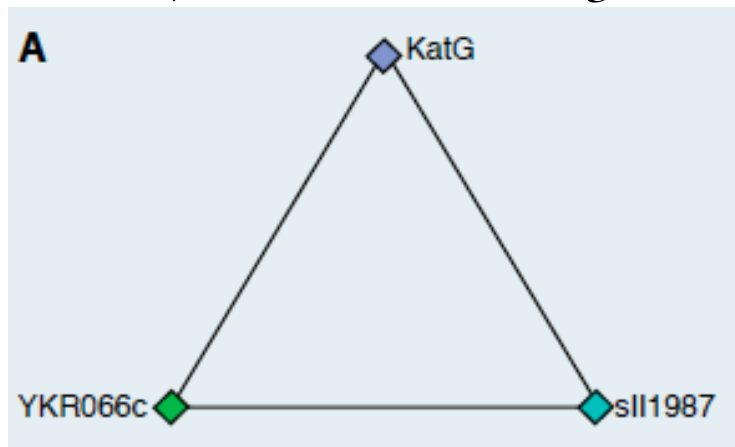


Boffelli et al. *Science* 299 :1391 (2003)

FAMILLES D'ORTHOLOGUES

BeT : [BLAST] *Best hit* (chaque “gène” de génome A comparé à ceux de génome B)

COGs (*Clusters of Orthologous Groups*) : déterminé par le graphe de BeTs

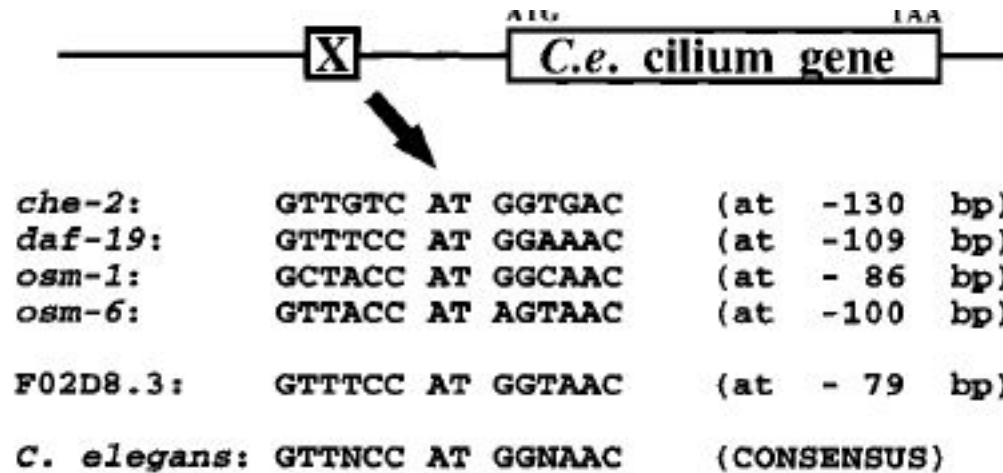


triangles formés par BeTs symétriques → squelette d'un COG
+ ajouter d'autres gènes avec des BeTs symétriques avec le groupe
+ inspection humaine

Tatusov & al. *Science* 278 :631 (1997)

RECHERCHE DE MOTIFS

Site de liaison de facteur de transcription : signal de régulation de l'expression génique



Placement inconnu, copies non-identiques

RECHERCHE DE MOTIFS 2

Problème abstrait : séquences S_1, \dots, S_n , avec des copies d'un motif M

Si on a une représentation du motif (signature, HMM, PSSM, ...) on sait comment trouver les copies

Qu'est-ce qu'on fait si on connaît pas M ?

Problème difficile...

Solution 1 : force brut — essayer tous les m -mers comme candidats de M , trouver les m -mers les plus proches à M en chaque séquence, et calculer un score (p.e., moyenne de distance de M , ou entropie)

Branch & bound

RECHERCHE DE MOTIFS 3

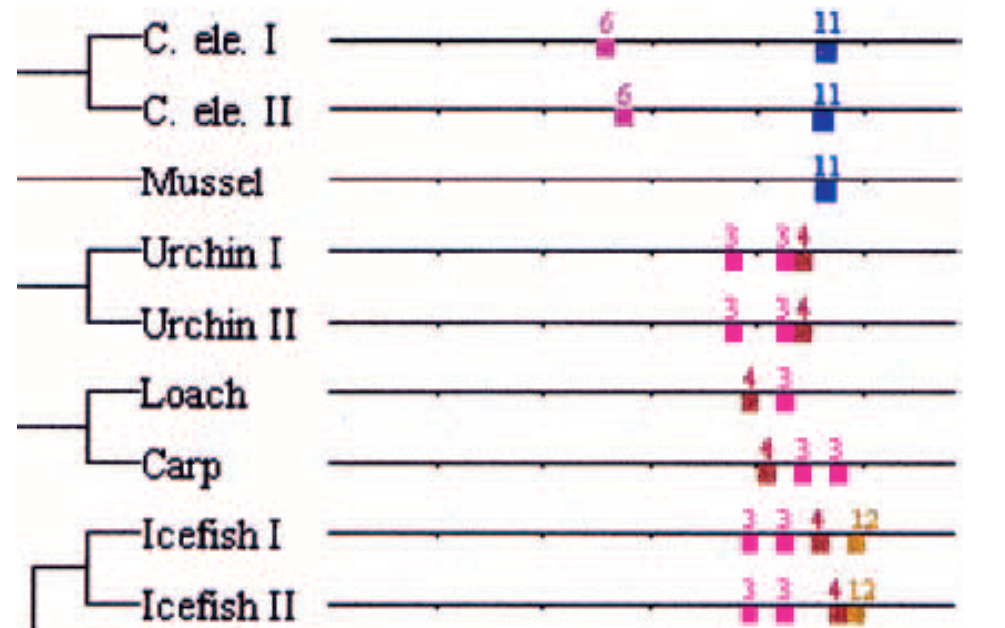
Score : étant donné des copies M_1, \dots, M_n , former profile et calculer le score comme LODS en utilisant un modèle de background

Solution 2 : Gibbs sampler

1. copies M_1, \dots, M_n
2. enlever M_j et le remplacer par le meilleur m -mer en S_j (comparaison au profile défini par $M_1, \dots, M_{j-1}, M_{j+1}, \dots, M_n$)
3. répéter jusqu'à convergence

PHYLOGENETIC FOOTPRINTING

Régions non-codantes conservées parmi des espèces distantes : éléments de régulation



Blanchette & Tompa, *Genome Res.* 12 : 739 (2002)