# **CONSERVATION DE SÉQUENCE**

Conservation \* IFT6299 H2014 \* UdeM \* Miklós Csűrös

#### Génomique comparative



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

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

#### L'ombre de sélection dans une fenêtre



hypothèse 0 = substitutions par  $e^{-\mu \mathbf{Q}t}$ hypothèse 1 = substitutions par  $e^{-\rho\mu\mathbf{Q}t}$ . Vraisemblances  $L_0, L_1$  et

$$LR = \log \frac{L_1}{L_0}$$

En haut : LLR dans un fenêtre de 50 pb glissant dans l'alignement multiple de 5 primates, avec des exons indiqués. En bas : % identité humain-souris.

Boffelli & al. Science 299 :1391 (2003)

## Puissance statistique

Détection dans une fenêtre de 50 pb selon % identité



Stone, Cooper & Sidow. Annu. Rev. Genomics Hum Genet. 6 :143, (2005).

# Puissance statistique II

Détection dans une fenêtre de 100 pb selon % identité



Stone, Cooper & Sidow. Annu Rev Genomics Hum Genet 6 :143, 2005).



Figure 1. Number of Genomes Required for Single Nucleotide Resolution

Eddy. PLoS Biology 3 :e10 (2005)



#### Conservation \* IFT6299 H2014 \* UdeM \* Miklós Csűrös

## Inférer la conservation

**Phylogenetic Tree** 

#### **Multiple Sequence Alignment**



Questions : (1) classes de taux fixées ou non; (2) estimation du taux neutre; (3) groupage de positions consécutives; (4) statistique pour détecter la conservation;
(5) alignement / trous / séquences manquantes

Davydov & al. PLoS Comput Biol 6 :e1001025 (2010)

# Contraintes fonctionnelles dans le génome humaine

Selon la méthode d'inférence et les données, on arrive à des conclusions différentes sur la fraction  $\alpha$  de régions conservées...



Estimates of  $\alpha_{sel}$  from 16 studies ranked by increasing values. Lower and upper bound values are indicated in blue and red, respectively.

Ponting & Hardison Genome Res 21 :1769 (2011)

# La méthode phylogénétique

On a un alignement multiple...



colonne = nucléotides homologues

 $\Rightarrow$  inférer la variation de contraintes évolutives dans le génome



- colonnes de résidues **homologues** \_\_\_\_\_

# Substitutions sur la phylogénie

on veut capturer la variation du taux de substitution  $\Rightarrow$  employer un modèle phylogénétique d'évolution

procéssus de Markov à temps continu avec matrice de taux  $\mathbf{Q}$ :

$$\mathbf{M}(t) = \begin{bmatrix} p_{A \to A}(t) & p_{A \to C}(t) & p_{A \to C}(t) & p_{A \to G}(t) \\ p_{C \to A}(t) & p_{C \to C}(t) & p_{C \to C}(t) & p_{C \to G}(t) \\ p_{G \to A}(t) & p_{G \to C}(t) & p_{G \to C}(t) & p_{G \to G}(t) \\ p_{T \to A}(t) & p_{T \to C}(t) & p_{T \to C}(t) & p_{T \to G}(t) \end{bmatrix} = \exp(\mathbf{Q}t)$$

 $\mathbf{M}(t)$  décrit les probabilités de substitution sur une arête de longueur t



### Taux variable — GERP++

(1) arbre neutre : longueur vient de sites «dégénerés» (p.e., codon GT N =valine) (2) définir la vraisemblance L(r) pour facteur d'échelle r > 0 — appliquer M(rt) sur arête de longueur t

(3) maximiser la fonction L(r) pour choisir r; score  $RS = (1-r) \sum_{uv \in ar \in tes} t_{uv}$ 



Davydov & al. PLoS Comput Biol 6 :e1001025 (2010)

## Distribution d'éléments conservés

Beaucoup d'éléments dans les régions introniques et intergéniques : régulation, gènes ARN



**B.** Composition of Constrained Elements

Davydov & al. PLoS Comput Biol 6 :e1001025 (2010)

# **Composition variable — SiPhy**

Utiliser une matrice de taux, varier échelle  $\omega$  (comme r avant) et distribution stationnaire  $\pi$  score dans une fenêtre de longueur k

$$\mathsf{LO} = \log \frac{\max_{\pi,\omega} \mathbb{P}\left\{ x_1[j..j+k-1], x_2[j..j+k-1], \dots, x_n[j..j+k-1] \mid \pi, \omega \right\}}{\mathbb{P}\left\{ x_1[j..j+k-1], x_2[j..j+k-1], \dots, x_n[j..j+k-1] \mid \pi_0, \omega_0 \right\}}$$



meilleure séparation par composition que par vitesse (4D = 3e position du codon dégéneré avec 4 choix; 2D = 3e position du codon dégéneré avec 2 choix)

Garber & al. Bioinformatics 25 :i64 (ISMB 2009)

## Effet de voisinage — SCONE

On veut capturer les dépendences entre nucléotides consécutifs :

(1) considérer des triples de nucléotides consécutifs, calculer les taux neutres pour  $AAA \rightarrow AAA$ ,  $AAA \rightarrow AAT$ , ...

(2) dans colonne j: reconstruire les caractères ancestraux par parcimonie dans colonnes j - 1 et j + 1, utiliser probabilité de substitutions

$$p_{x[j] \to y[j]} = \frac{\mathbb{P}\left\{x[j-1..j+1] \to y[j-1..j+1]\right\}}{\sum_{a,b} \mathbb{P}\left\{x[j-1..j+1] \to ay[j]b\right\}}$$

score = espérance postérieure avec m catégories d'échelle  $\omega_k$  (pondération par vraisemblance  $L(\omega_k)$ )

$$\mathbb{E}\omega = \frac{\sum_{k=1}^{m} \omega_k \cdot L(\omega_k)}{\sum_{k=1}^{m} L(\omega_k)}$$

P-value : vérifier distribution de  $\mathbb{E}\omega$  dans évolution simulée (10000 fois) sur l'arbre neutre

#### **SCONE**

#### (intron et une région intergénique)



Asthana & al. PLoS Computational Biology 3 :e254 (2007)

# Effet de voisinage — PhyloHMM

phastCons



émissions : colonnes de l'alignement multiple, avec probabilités de transition  $e^{\mathbf{Q}t}$ (neutre) ou  $e^{\rho \mathbf{Q}t}$  (sélection négative avec  $\rho < 0$ )

Siepel & al Genome Res 15 :1034 (2005)

#### Turnover

Modélisation par phylo-HMM



Siepel & al. RECOMB 2006



Siepel & al. RECOMB 2006

## Accélération — HAR

Human Accelerated Regions : plus rapide dans l'humain mais conservé ailleurs

Vol 443|14 September 2006|doi:10.1038/nature05113

#### ARTICLES

nature

# An RNA gene expressed during cortical development evolved rapidly in humans

Katherine S. Pollard<sup>1</sup>\*†, Sofie R. Salama<sup>1,2</sup>\*, Nelle Lambert<sup>4,5</sup>, Marie-Alexandra Lambot<sup>4</sup>, Sandra Coppens<sup>4</sup>, Jakob S. Pedersen<sup>1</sup>, Sol Katzman<sup>1</sup>, Bryan King<sup>1,2</sup>, Courtney Onodera<sup>1</sup>, Adam Siepel<sup>1</sup>†, Andrew D. Kern<sup>1</sup>, Colette Dehay<sup>6,7</sup>, Haller Igel<sup>3</sup>, Manuel Ares Jr<sup>3</sup>, Pierre Vanderhaeghen<sup>4</sup> & David Haussler<sup>1,2</sup>



Amadio & Walsh Cell 126 :1033 (2006)

# Pas trop grande différence avec assez de génomes ...



**Figure 1.** Receiver operating characteristic (ROC) curves showing falsepositive versus true-positive rates for the all-branch tests implemented in phyloP: (red) LRT, (green) SCORE, (blue) SPH, and (purple) GERP. Individual plots show results for simulated data sets with either 3-bp (*top*) or 1-bp (*bottom*) elements generated from models with a range of deviations  $\rho$  from the neutral rate  $\rho = 1.0$  (columns).

Pollard & al. Genome Res 20 :110 (2010)

# Accord qualitatif sur sélection

#### **Table 1.** Publications describing estimates of $\alpha_{sel}$ and $\alpha_{sel}^0$

Publication	Method ( $\alpha_{sel}$ and $\alpha_{sel}^{0}$ , estimation)	α <sub>sel</sub> (%) Iower	α <sub>sel</sub> (%) higher	Substitutions or indels or topography	Neutral model/standard	Whole or partial genome	Multiple, or pair of, genomes	Local or global neutral rate
Lunter et al. (2006)	NIM ( $\alpha_{sel}$ )	2.56	3.25	Indels	Randomly placed indels	Whole	Pair	Local
Thomas et al. (2003)	MCSs ( $\alpha_{sel}$ )	3.7	3.7	Substitutions	4D sites	Partial	Multiple	Local
The ENCODE Project Consortium (2007) (ENCODE)	Two of three methods ( $\alpha_{sel}$ )	4.9	4.9	Substitutions	4D sites/most aligned sites	Partial (ENCODE)	Multiple	Global
Lindblad-Toh et al. (2005)	Substitutions ( $\alpha_{sel}$ )	5.3	5.3	Substitutions	ARs	Whole	Pair	Local
Pollard et al. (2010)	Various $(\alpha_{sel})$	5.3	5.3	Substitutions	Various	Partial (ENCODE)	Multiple	Global
Lindblad-Toh et al. (2011)	SiPhy ( $\alpha_{sel}$ )	5.4	5.4	Patterns	ARs	Whole	Multiple	Global
Eory et al. (2010)	Substitutions ( $\alpha_{sel}$ )	5.4	5.4	Substitutions	ARs	Whole	Pair	Local
Cooper et al. (2005)	GERP ( $\alpha_{sel}$ )	5.5	5.5	Substitutions	Most aligned sites	Partial	Multiple	Local
Chiaromonte et al. (2003)	Substitutions ( $\alpha_{sel}$ )	2.29	6.15	Substitutions	ARs	Whole	Pair	Local
Siepel et al. (2005)	phastCons ( $\alpha_{sel}$ )	3	8	Substitutions	none	Whole	Multiple	Global
Davydov et al. (2010)	GERP++ ( $\alpha_{sel}$ )	6	8	Substitutions	Most aligned sites	Whole	Multiple	Global
Smith et al. (2004)	Substitutions/Turnover $(\alpha_{sel} \text{ and } \alpha_{sel}^{0})$	10	10	Substitutions	Simulation	Partial	Multiple pairs	Local
Meader et al. (2010)	NIM ( $\alpha_{sel}$ and $\alpha^{0}_{sel}$ )	6.5	10	Indels	Randomly placed indels	Whole	Multiple pairs	Local
Garber et al. (2009)	SiPhy ( $\alpha_{sel}$ )	5.8	10.2	Patterns	ARs	Partial (ENCODE)	Multiple	Global
Asthana et al. (2007)	SCONE ( $\alpha_{sel}$ )	5.5	11	Substitutions within trinucleotides	Most aligned sites	Partial (ENCODE)	Multiple	Global
Parker et al. (2009)	Topography (Chai; $\alpha_{sel}$ )	12	12	Topography	Most aligned sites	Partial (ENCODE)	Multiple	Global

(4D) Four-fold degenerate; (ARs) ancestral repeats; (ENCODE) Encyclopedia of DNA Elements project; (GERP) Genomic Evolutionary Rate Profiling; (MCSs) multispecies conserved sequence; (NIM) neutral indel model; (SCONE) Sequence CONservation Evaluation.

Ponting & Hardison Genome Res 21 :1769 (2011)